

Characterizing Truthful Multi-Armed Bandit Mechanisms*

Moshe Babaioff
Microsoft Research
Mountain View, CA 94043
moshe@microsoft.com

Yogeshwer Sharma[†]
Cornell University
Ithaca, NY 14853
yogi@cs.cornell.edu

Aleksandrs Slivkins
Microsoft Research
Mountain View, CA 94043
slivkins@microsoft.com

November 2008

Minor revisions: February 2009, June 2009, Sept 2009

Abstract

We consider a multi-round auction setting motivated by pay-per-click auctions for Internet advertising. In each round the auctioneer selects an advertiser and shows her ad, which is then either clicked or not. An advertiser derives value from clicks; the value of a click is her private information. Initially, neither the auctioneer nor the advertisers have any information about the likelihood of clicks on the advertisements. The auctioneer’s goal is to design a (dominant strategies) truthful mechanism that (approximately) maximizes the social welfare.

If the advertisers bid their true private values, our problem is equivalent to the *multi-armed bandit problem*, and thus can be viewed as a strategic version of the latter. In particular, for both problems the quality of an algorithm can be characterized by *regret*, the difference in social welfare between the algorithm and the benchmark which always selects the same “best” advertisement. We investigate how the design of multi-armed bandit algorithms is affected by the restriction that the resulting mechanism must be truthful. We find that truthful mechanisms have certain strong structural properties – essentially, they must separate exploration from exploitation – *and* they incur much higher regret than the optimal multi-armed bandit algorithms. Moreover, we provide a truthful mechanism which (essentially) matches our lower bound on regret.

ACM Categories and subject descriptors: F.2.2 [Analysis of Algorithms and Problem Complexity]: Nonnumerical Algorithms and Problems; K.4.4 [Computers and Society]: Electronic Commerce; F.1.2 [Computation by Abstract Devices]: Modes of Computation—*Online computation*; J.4 [Social and Behavioral Sciences]: Economics

General Terms: theory, algorithms, economics.

Keywords: mechanism design, truthful mechanisms, single-parameter auctions, multi-armed bandit problem, regret, online learning.

*A preliminary version [9] of this paper has appeared in ACM EC 2009.

[†]This research was done while the author was an intern at Microsoft Research, Silicon Valley Center.

1 Introduction

In recent years there has been much interest in understanding the implication of strategic behavior on the performance of algorithms whose input is distributed among selfish agents. This study was mainly motivated by the Internet, the main arena of large scale interaction of agents with conflicting goals. The field of Algorithmic Mechanism Design [35] studies the design of mechanisms in computational settings (for background see the recent book [36] and survey [38]).

Much attention has been drawn to the market for sponsored search (e.g. [28, 19, 39, 32, 3]), a billions dollar market with numerous auctions running every second. Research on sponsored search mostly focus on equilibria of the Generalized Second Price (GSP) auction [19, 39], the auction that is most commonly used in practice (e.g. by Google and Yahoo), or on the design of truthful auctions [2]. All these auctions rely on knowing the rates at which users click on the different advertisements (a.k.a. Click-Through-Rates, or CTRs), and do not consider the process in which these CTRs are learned or refined over time by observing users' behavior. We argue that strategic agents would take this process into account, as it influences their utility. Prior work [22] focused on the implication of click fraud on the methods used to learn CTRs. We on the other hand are interested in the implications of the *strategic bidding* by the agents. Thus, we consider the problem of designing truthful sponsored search auctions when the process of learning the CTRs is a part of the game.

We are mainly interested in the interplay between the online learning and the strategic aspects of the problem. To isolate this issue, we consider the following setting, which is a natural *strategic* version of the multi-armed bandit (MAB) problem. In this setting, there are k agents. Each agent i has a single advertisement, and a *private* value $v_i > 0$ for every click she gets. The mechanism is an online algorithm that first solicits bids from the agents, and then runs for T rounds. In each round the mechanism picks an agent (using the bids and the clicks observed in the past rounds), displays her advertisement, and receives a feedback – if there was a click or not. Payments are assigned after round T . Each agent tries to maximize her own utility: the difference between the value that she derives from clicks and the payment she pays. We assume that initially no information is known about the likelihood of each agent to be clicked, and in particular there are no Bayesian priors.

We are interested in designing mechanisms which are truthful (in dominant strategies): every agent maximizes her utility by bidding truthfully, for any bids of the others and *for any clicks* that would have been received. The goal is to maximize the social welfare.¹ Since the payments cancel out, this is equivalent to maximizing the total value derived from clicks, where an agent's contribution to that total is her private value times the number of clicks she receives. We call this setting the *MAB mechanism design problem*.

In the absence of strategic behavior this problem reduces to a standard MAB formulation in which an algorithm repeatedly chooses one of the k alternatives (“arms”) and observes the associated payoff: the value-per-click of the corresponding ad if the ad is clicked, and 0 otherwise. The crucial aspect in MAB problems is the tradeoff between acquiring more information (*exploration*) and using the current information to choose a good agent (*exploitation*). MAB problems have been studied intensively for the past three decades (see [13, 14, 20]). In particular, the above formulation is well-understood [6, 7, 16] in terms of *regret* relative to the benchmark which always chooses the same “best” alternative. This notion of regret naturally extends to the strategic setting outlined above, the total payoff being exactly equal to the social welfare, and the regret being exactly the loss in social welfare. Thus one can directly compare MAB algorithms and MAB mechanisms in terms of welfare loss (regret).

Broadly, we ask how the design of MAB algorithms is affected by the restriction of truthfulness: what is the difference between the best *algorithms* and the best *truthful mechanisms*? We are interested both in

¹Social welfare includes both the auctioneer's revenue and the agents' utility. Since in practice different sponsored search platforms compete against one another, taking into account the agents' utility increases the platform's attractiveness to the advertisers.

terms of the structural properties and the gap in performance (in terms of regret). We are not aware of any prior work that characterizes truthful learning algorithms or proves negative results on their performance.

Our contributions. We present two main contributions. First, we present a characterization of (dominant-strategy) truthful mechanisms. Second, we present a lower bound on the regret that such mechanisms must suffer. This regret is significantly larger than the regret of the best MAB algorithms.

Formally, a mechanism for the MAB mechanism design problem is a pair $(\mathcal{A}, \mathcal{P})$, where \mathcal{A} is the *allocation rule* (essentially, an MAB algorithm), and \mathcal{P} is the *payment rule*. Note that regret is completely determined by the allocation rule. As is standard in the literature, we focus on mechanisms in which each agent’s payment (averaged over clicks) is between 0 and her bid; such mechanisms are called *normalized*, and they satisfy voluntary participation.

The setting we study is a *single-parameter auction*, the most studied and well-understood type of auctions. For such settings truthful mechanisms are fully characterized [33, 4]: a mechanism is truthful if and only if the allocation rule is monotone (by increasing her bid an agent cannot cause a decrease in the number of clicks she gets), and the payment rule is defined in a specific and (essentially) unique way. Yet, this characterization is *not* the right characterization for the MAB setting! The main problem is that in our setting click information for any agent that is not chosen at a given round is not available to the mechanism, and thus cannot be used in the computation of payments. Thus, the payment cannot depend on any unobserved clicks. We show that this has severe implications on the structure of truthful mechanisms.

The first notable property of a truthful mechanism is a much stronger version of monotonicity:

Definition 1.1. A *realization* consists of click information for all agents at all rounds (including unobserved ones). An allocation rule is *pointwise monotone* if for each realization, each bid profile and each round, if an agent is played at the round, then she is also played after increasing her bid (fixing everything else).

Let us consider (for the ease of exposition) allocation rules that satisfy the following two natural conditions. First, an allocation rule is *scale-free* if it is invariant under multiplying all bids by the same positive number (essentially, changing the currency unit). Second, it is *Independent of Irrelevant Alternatives (IIA)*, for short) if for any given realization, bid profile and round, a change of bid of agent i cannot transfer the allocation in this round from agent j to agent l , where these are three distinct agents.

We show that any truthful mechanism must have a strict separation between exploration and exploitation. A crucial feature of exploration is the ability to influence the allocation in forthcoming rounds. To make this point more concrete, we call a round *influential* for a given realization if for some bid profile changing the realization for this round can affect the allocation in some future round. We show that in any such round, the allocation can not depend on the bids. Thus, influential rounds are essentially useless for exploitation.

Definition 1.2. An allocation rule \mathcal{A} is called *exploration-separated* if for any given realization, the allocation in any influential round for that realization does not depend on the bids.

We are now ready to present our main structural result, which is in fact a complete characterization.

Theorem 1.3. Consider the MAB mechanism design problem. Let \mathcal{A} be a non-degenerate² deterministic allocation rule which is scale-free and satisfies IIA. Then mechanism $(\mathcal{A}, \mathcal{P})$ is normalized and truthful for some payment rule \mathcal{P} if and only if \mathcal{A} is pointwise monotone and exploration-separated.

²Non-degeneracy is a mild technical assumption, formally defined in “preliminaries”, which ensures that (essentially) if a given allocation happens for some bid profile (b_i, b_{-i}) then the same allocation happens for all bid profiles (x, b_{-i}) , where x ranges over some non-degenerate interval. Without this assumption, all structural results hold (essentially) *almost surely* w.r.t the k -dimensional Lebesgue measure on the bid vectors. Exposition becomes significantly more cumbersome, yet leads to the same lower bounds on regret. For clarity, we assume non-degeneracy throughout this version of the paper.

We also obtain a similar (but somewhat more complicated) characterization without assuming that allocations are scale-free and satisfy IIA (Theorem 3.8). We use it then to derive Theorem 1.3. We emphasize that our characterization results hold regardless of whether the auctioneer’s goal is to maximize welfare or revenue or any other objective.

In view of Theorem 1.3, we present a lower bound on the performance of exploration-separated algorithms. We consider a setting, termed the *stochastic MAB mechanism design problem*, in which each click on a given advertisement is an independent random event which happens with a fixed probability, a.k.a. the CTR. The expected “payoff” from choosing a given agent is her private value times her CTR. For the ease of exposition, assume that the bids lie in the interval $[0, 1]$. Then the non-strategic version is the *stochastic MAB problem* in which the payoff from choosing a given arm i is an independent sample in $[0, 1]$ with a fixed mean μ_i . In both versions, *regret* is defined with respect to a hypothetical allocation rule (resp. algorithm) that always chooses an arm with the maximal expected payoff. Specifically, regret is the expected difference between the social welfare (resp. total payoff) of the benchmark and that of the allocation rule (resp. algorithm). The goal is to minimize $R(T)$, worst-case regret over all problem instances on T rounds.

We show that the worst-case regret of any exploration-separated mechanism is *larger* than that of the optimal MAB algorithm [7]: $\Omega(T^{2/3})$ vs $O(\sqrt{T})$ for a fixed number of agents. We obtain an even more pronounced difference if we restrict our attention to the δ -gap problem instances: instances for which the best agent is better than the second-best by a (comparatively large) amount δ , that is $\mu_1 v_1 - \mu_2 v_2 = \delta \cdot (\max_i v_i)$, where arms are arranged such that $\mu_1 v_1 \geq \mu_2 v_2 \geq \dots \geq \mu_k v_k$. Such instances are known to be easy for the MAB algorithms. Namely, an algorithm can achieve the optimal worst-case regret $O(\sqrt{kT \log T})$ and regret $O(\frac{k}{\delta} \log T)$ on δ -gap instances [29, 6]. However, for exploration-separated mechanisms the worst-case regret $R_\delta(T)$ over the δ -gap instances is polynomial in T as long as worst-case regret is even remotely non-trivial (i.e., sublinear). Thus, for the δ -gap instances the gap between algorithms and truthful mechanisms in the worst-case regret is *exponential* in T .

Theorem 1.4. *Consider the stochastic MAB mechanism design problem with k agents. Let \mathcal{A} be a deterministic allocation rule that is exploration-separated. Then \mathcal{A} has worst-case regret $R(T) = \Omega(k^{1/3} T^{2/3})$. Moreover, if $R(T) = O(T^\gamma)$ for some $\gamma < 1$ then for every fixed $\delta \leq \frac{1}{4}$ and $\lambda < 2(1 - \gamma)$ the worst-case regret over the δ -gap instances is $R_\delta(T) = \Omega(\delta T^\lambda)$.*

We note that our lower bounds holds for a more general setting in which the values-per-click can change over time, and the advertisers are allowed to change their bids at every time step.

To complete the picture, we present a very simple (deterministic) mechanism that is truthful and normalized, and matches the lower bound $R(T) = \Omega(k^{1/3} T^{2/3})$ up to logarithmic factors.

We also provide a number of extensions. First, we prove a similar (but slightly weaker) regret bound without the scale-free assumption. Second, we extend some of our results to randomized mechanisms; in this setting, (dominant-strategy) truthfulness means “truthfulness for each realization of the private randomness”. Third, we consider a weaker notion of truthfulness for randomized mechanisms – for each realization of the clicks, but in expectation over the random seed, and use this notion to provide algorithmic results for the version of the MAB mechanism design problem in which clicks are chosen by an adversary. Fourth, we discuss an even more permissive notion of truthfulness – truthfulness in expectation over the clicks (and the random seed).

Other related work and discussion. The question of how the performance of a truthful mechanism compares to that of the optimal algorithm for the corresponding non-strategic problem has been considered in the literature in a number of other auction settings. Performance gaps have been shown for various scheduling problems [4, 35, 18] and for online auction for expiring goods [31]. Other papers presented approximation gaps due to *computational constraints*, e.g. for combinatorial auctions [30, 18] and combinatorial public projects [37], showing a gap via a structural result for truthful mechanisms.

The study of MAB mechanisms has been initiated by Gonen and Pavlov [21]. The authors present a MAB mechanism which is claimed to be truthful in a certain approximate sense. Unfortunately, this mechanism does not satisfy the claimed properties; this was also confirmed with the authors through personal communication (see also a similar note in [17]).

MAB algorithms were used in the design of Cost-Per-Action sponsored search auctions in Nazerzadeh et al. [34], where the authors construct a mechanism with approximate properties of truthfulness and individual rationality. Approximately truthful mechanisms are reasonable assuming the agents would not lie unless it leads to significant gains. However, this solution concept is weaker than the exact notion and it may still be rational for the agents to deviate (perhaps significantly) from being truthful. Moreover, as truthful bidding is not a Nash equilibrium, agents might have an increased incentive to deviate if they speculate that others are deviating. All of that may result in unpredictable, and possibly highly suboptimal outcomes. In this paper we focus on understanding what can be achieved with the *exact* truthfulness, mainly proving results of structural and lower-bounding nature. We note in passing that providing similar results for the approximately truthful setting such as the one in [34] is a worthy and challenging open question.

Independently and concurrently, Devanur and Kakade [17] have studied truthful MAB mechanisms with focus on maximizing the revenue. They present a lower bound of $\Omega(T^{2/3})$ on the loss in revenue with respect to the VCG (Vickrey-Clarke-Groves) payment, as well as a truthful mechanism that matches the lower bound. (This mechanism is almost identical to the one that we present in order to match the lower bound in Theorem 4.1.)

Our lower bounds use (a novel application of) the relative entropy technique from [29, 7], see [26] for an account. For other application of this technique, see e.g. [16, 23, 27, 11].

Our work focuses on regret in a prior-free setting in which the algorithm has no prior on CTRs. This is in contrast to the recent line of work on *dynamic auctions* [12, 5] which considers fully Bayesian settings in which there is a known prior on CTRs, and VCG-like social welfare-maximizing mechanisms are feasible. In our prior-free setting VCG-mechanisms cannot be applied as such mechanisms require the allocation to exactly maximize the expected social welfare, which is impossible (and not well-defined) without a prior.

We require the mechanisms to satisfy a strong notion of truthfulness: bidding truthful is optimal for *every* possible realization (and bids of others). This notion is attractive as it does not require the agents to be risk neutral. Moreover, it allows for the CTRs to change over time (and still incentivizes agents to be truthful). Finally, an agent never regrets in retrospect that she has been truthful. It is desirable to understand this notion before moving to weaker notions.

Map of the paper. Section 2 is preliminaries. Truthfulness characterization is developed and proved in Section 3. The lower bounds on regret and the simple mechanism that matches them are in Section 4. Extensions and open questions are in Section 5. To improve the flow of the paper, some of the material is moved to the appendices.

2 Definitions and preliminaries

In the MAB mechanism design problem, there is a set K of k agents numbered from 1 to k . Each agent i has a *value* $v_i > 0$ for every click she gets; this value is known only to agent i . Initially, each agent i submits a *bid* $b_i > 0$, possibly different from v_i .³ The “game” lasts for T rounds, where T is the given *time horizon*. A *realization* represents the click information for all agents and all rounds. Formally, it is a tuple

³One can also consider a more realistic and general model in which the value-per-click of an agent changes over time and the agents are allowed to change their bid at every round. The case that the value-per-click of each agent does not change over time is a special case. In that case truthfulness implies that each agent basically submits one bid as in our model (the same bid at every round), thus our main results (necessary conditions for truthfulness and regret lower bounds) also hold for the more general model.

$\rho = (\rho_1, \dots, \rho_k)$ such that for every agent i and round t , the bit $\rho_i(t) \in \{0, 1\}$ indicates whether i gets a click if played at round t . An *instance* of the MAB mechanism design problem consists of the number of agents k , time horizon T , a vector of private values $v = (v_1, \dots, v_k)$, a vector of bids (*bid profile*) $b = (b_1, \dots, b_k)$, and realization ρ .

A *mechanism* is a pair $(\mathcal{A}, \mathcal{P})$, where \mathcal{A} is allocation rule and \mathcal{P} is the payment rule. An *allocation rule* is represented by a function \mathcal{A} that maps bid profile b , realization ρ and a round t to the agent i that is chosen (receives an *impression*) in this round: $\mathcal{A}(b; \rho; t) = i$. We also denote $\mathcal{A}_i(b; \rho; t) = \mathbf{1}_{\{\mathcal{A}(b; \rho; t) = i\}}$. The allocation is *online* in the sense that at each round it can only depend on clicks observed prior to that round. Moreover, it does not know the realization in advance; in every round it only observes the realization for the agent that is shown in that round. A *payment rule* is a tuple $\mathcal{P} = (\mathcal{P}_1, \dots, \mathcal{P}_k)$, where $\mathcal{P}_i(b; \rho) \in \mathbb{R}$ denotes the payment charged to agent i when the bids are b and the realization is ρ .⁴ Again, the payment can only depends on observed clicks. A mechanism is called *normalized* if for any agent i , bids b and realization ρ it holds that $\mathcal{P}_i(b; \rho)$ is non-negative and at most b_i times the number of clicks agent i got.

For given realization ρ and bid profile b , the number of clicks received by agent i is denoted $\mathcal{C}_i(b; \rho)$. Call $\mathcal{C} = (\mathcal{C}_1, \dots, \mathcal{C}_k)$ the *click-allocation* for \mathcal{A} . The *utility* that agent i with value v_i gets from the mechanism $(\mathcal{A}, \mathcal{P})$ when the bids are b and the realization is ρ is $\mathcal{U}_i(v_i; b; \rho) = v_i \cdot \mathcal{C}_i(b; \rho) - \mathcal{P}_i(b; \rho)$ (quasi-linear utility). The mechanism is *truthful* if for any agent i , value v_i , bid profile b and realization ρ it is the case that $\mathcal{U}_i(v_i; v_i, b_{-i}; \rho) \geq \mathcal{U}_i(v_i; b_i, b_{-i}; \rho)$.

In the *stochastic* MAB mechanism design problem, an adversary specifies a vector $\mu = (\mu_1, \dots, \mu_k)$ of CTRs (concealed from \mathcal{A}), then for each agent i and round t , realization $\rho_i(t)$ is chosen independently with mean μ_i . Thus, an instance of the problem includes μ rather than a fixed realization. For a given problem instance \mathcal{I} , let $i^* \in \arg\max_i \mu_i v_i$, then *regret* on this instance is defined as

$$R^{\mathcal{I}}(T) = T v_{i^*} \mu_{i^*} - \mathbb{E} \left[\sum_{t=1}^T \sum_{i=1}^k \mu_i v_i \mathcal{A}_i(b; \rho; t) \right]. \quad (2.1)$$

For a given parameter v_{\max} , the *worst-case regret*⁵ $R(T; v_{\max})$ denotes the supremum of $R^{\mathcal{I}}(T)$ over all problem instances \mathcal{I} in which all private values are at most v_{\max} . Similarly, we define $R_{\delta}(T; v_{\max})$, the *worst-case δ -regret*, by taking the supremum only on instances with δ -gap.

Most of our results are stated for *non-degenerate* allocation rules, defined as follows. An interval is called *non-degenerate* if it has positive length. Fix bid profile b , realization ρ , and rounds t and t' with $t \leq t'$. Let $i = \mathcal{A}(b; \rho; t)$ and ρ' be the allocation obtained from ρ by flipping the bit $\rho_i(t)$. An allocation rule \mathcal{A} is *non-degenerate* w.r.t. (b, ρ, t, t') if there exists a non-degenerate interval I containing b_i such that

$$\mathcal{A}_i(x, b_{-i}; \varphi; s) = \mathcal{A}_i(b; \varphi; s) \quad \text{for each } \varphi \in \{\rho, \rho'\}, \text{ each } s \in \{t, t'\}, \text{ and all } x \in I.$$

An allocation rule is *non-degenerate* if it is non-degenerate w.r.t. each tuple (b, ρ, t, t') .

3 Truthfulness characterization

Before presenting our characterization we begin by describing some related background. The click allocation \mathcal{C} is *non-decreasing* if for each agent i , increasing her bid (and keeping everything else fixed) does not decrease \mathcal{C}_i . Prior work has established a characterization of truthful mechanisms for single-parameter domains (domains in which the private information of each agent is one-dimensional), relating click allocation monotonicity and truthfulness (see below). For our problem, this result is a characterization of MAB

⁴We allow the mechanism to determine the payments at the end of the T rounds, and not after every round. This makes that task of designing a truthful mechanism *easier* and thus strengthen our necessary condition for truthfulness (the condition used to derive the lower bounds on regret.)

⁵By abuse of notation, when clear from the context, the “worst-case regret” is sometimes simply called “regret”.

algorithms that are truthful for a given realization ρ , assuming that the *entire* realization ρ can be used to compute payments (when computing payments one can use click information for every round and every agent, even if the agent was not shown at that round.) One of our main contributions is a characterization of MAB allocation rules that can be truthfully implemented when payment computation is restricted to only use clicks information of the actual impressions assigned by the allocation rule.

An MAB allocation rule \mathcal{A} is *truthful with unrestricted payment computation* if it is truthful with a payment rule that can use the *entire* realization ρ in its computation. We next present the prior result characterizing truthful mechanisms with unrestricted payment computation.

Theorem 3.1 (Myerson [33], Archer and Tardos [4]). *Let $(\mathcal{A}, \mathcal{P})$ be a normalized mechanism for the MAB mechanism design problem. It is truthful with unrestricted payment computation if and only if for any given realization ρ the corresponding click-allocation \mathcal{C} is non-decreasing and the payment rule is given by*

$$\mathcal{P}_i(b_i, b_{-i}; \rho) = b_i \cdot \mathcal{C}_i(b_i, b_{-i}; \rho) - \int_0^{b_i} \mathcal{C}_i(x, b_{-i}; \rho) dx. \quad (3.1)$$

We can now move to characterize truthful MAB mechanisms when the payment computation is restricted. The following notation will be useful: for a given realization ρ , let $\rho \oplus \mathbf{1}(i, t)$, be the realization that coincides with ρ everywhere, except that the bit $\rho_i(t)$ is flipped.

The first notable property of truthful mechanisms is a stronger version of monotonicity. Recall (see Definition 1.1) that an allocation rule \mathcal{A} is *pointwise monotone* if for each realization ρ , bid profile b , round t and agent i , if $\mathcal{A}_i(b_i, b_{-i}; \rho; t) = 1$ then $\mathcal{A}_i(b_i^+, b_{-i}; \rho; t) = 1$ for any $b_i^+ > b_i$. In words, increasing a bid cannot cause a loss of an impression.

Lemma 3.2. *Consider the MAB mechanism design problem. Let $(\mathcal{A}, \mathcal{P})$ be a normalized truthful mechanism such that \mathcal{A} is a non-degenerate deterministic allocation rule. Then \mathcal{A} is pointwise-monotone.*

Proof. For a contradiction, assume not. Then there is a realization ρ , a bid profile b , a round t and agent i such that agent i loses an impression in round t by increasing her bid from b_i to some larger value b_i^+ . In other words, we have $\mathcal{A}_i(b_i^+, b_{-i}; \rho; t) < \mathcal{A}_i(b_i, b_{-i}; \rho; t)$. Without loss of generality, let us assume that there are no clicks after round t , that is $\rho_j(t') = 0$ for any agent j and any round $t' > t$ (since changes in ρ after round t does not affect anything before round t).

Let $\rho' = \rho \oplus \mathbf{1}(i, t)$. The allocation in round t cannot depend on this bit, so it must be the same for both realizations. Now, for each realization $\varphi \in \{\rho, \rho'\}$ the mechanism must be able to compute the price for agent i when bids are (b_i^+, b_{-i}) . That involves computing the integral $I_i(\varphi) = \int_{x \leq b_i^+} \mathcal{C}_i(x, b_{-i}; \varphi) dx$ from (3.1). We claim that $I_i(\rho) \neq I_i(\rho')$. However, the mechanism cannot distinguish between ρ and ρ' since they only differ in bit (i, t) and agent i does not get an impression in round t . This is a contradiction.

It remains to prove the claim. Without loss of generality, assume that $\rho_i(t) = 0$ (otherwise interchange the role of ρ and ρ'). We first note that $\mathcal{C}_i(x, b_{-i}; \rho) \leq \mathcal{C}_i(x, b_{-i}; \rho')$ for every x . This is because everything is same in ρ and ρ' until round t (so the impressions are same too), there are no clicks after round t , and in round t the behavior of \mathcal{A} on the two realizations can be different only if that agent i gets an impression, in which case she is clicked under ρ' and not clicked under ρ .

Since \mathcal{A} is non-degenerate, there exists a non-degenerate interval I containing b_i such that changing bid of agent i to any value in this interval does not change the allocation at round t (both for ρ and for ρ'). For any $x \in I$ we have $\mathcal{C}_i(x, b_{-i}; \rho) < \mathcal{C}_i(x, b_{-i}; \rho')$, where the difference is due to the click in round t . It follows that $I_i(\rho) < I_i(\rho')$. Claim proved. Hence, the mechanism cannot be implemented truthfully. \square

Recall (see Definition 1.2) that round t is *influential* for a given realization ρ if for some bid profile b there exists a round $t' > t$ such that $\mathcal{A}(b; \rho; t') \neq \mathcal{A}(b; \rho \oplus \mathbf{1}(j, t); t')$ for $j = \mathcal{A}(b; \rho; t)$. In words: changing the relevant part of the realization at round t affects the allocation in some future round t' . An allocation

rule \mathcal{A} is called *exploration-separated* if for any given realization ρ and round t that is influential for ρ , it holds that $\mathcal{A}(b; \rho; t) = \mathcal{A}(b'; \rho; t)$ for any two bid vectors b, b' (allocation at t does not depend on the bids).

The main structural implication is “truthful implies exploration-separated”. To illustrate the ideas behind this implication, we first state and prove it for two agents.

Proposition 3.3. *Consider the MAB mechanism design problem with two agents. Let \mathcal{A} be a non-degenerate scale-free deterministic allocation rule. If $(\mathcal{A}, \mathcal{P})$ is a normalized truthful mechanism for some \mathcal{P} , then it is exploration separated.*

Proof. Assume \mathcal{A} is not exploration-separated. Then there is a *counterexample* (ρ, t) : a realization ρ and a round t such that round t is influential and allocation in round t depends on bids. We want to prove that this leads to a contradiction.

Let us pick a counterexample (ρ, t) with some useful properties. Since round t is influential, there exists a realization ρ and bid profile b such that the allocation at some round $t' > t$ (the *influenced* round) is different under realization ρ and another realization $\rho' = \rho \oplus \mathbf{1}(j, t)$, where $j = \mathcal{A}(b; \rho; t)$ is the agent chosen at round t under ρ . Without loss of generality, let us pick a counterexample with minimum value of t' over all choices of (b, ρ, t) . For ease of exposition, from this point on let us assume that $j = 2$. For the counterexample we can also assume that $\rho_1(t') = 1$, and that there are no clicks after round t' , that is $\rho_l(t'') = \rho'_l(t'') = 0$ for all $t'' > t'$ and for all $l \in \{1, 2\}$.

We know that the allocation in round t depends on bids. This means that agent 1 gets an impression in round t for some bid profile $\hat{b} = (\hat{b}_1, \hat{b}_2)$ under realization ρ , that is $\mathcal{A}(\hat{b}; \rho; t) = 1$. As the mechanism is scale-free this means that, denoting $b_1^+ = \hat{b}_1 b_2 / \hat{b}_2$ we have $\mathcal{A}(b_1^+, b_2; \rho; t) = 1$. Since $\mathcal{A}(b_1, b_2; \rho; t) = 2$ and $\mathcal{A}(b_1^+, b_2; \rho; t) = 1$, pointwise monotonicity (Lemma 3.2) implies that $b_1^+ > b_1$. We conclude that there exists a bid $b_1^+ > b_1$ for agent 1 such that $\mathcal{A}(b_1^+, b_2; \rho; t) = 1$.

Now, the mechanism needs to compute prices for agent 1 for bids (b_1^+, b_2) under realizations ρ and ρ' , that is $\mathcal{P}_1(b_1^+, b_2; \rho)$ and $\mathcal{P}_1(b_1^+, b_2; \rho')$. Therefore, the mechanism needs to compute the integral $I_1(\varphi) = \int_{x \leq b_1^+} \mathcal{C}_1(x, b_2; \varphi) dx$ for both realizations $\varphi \in \{\rho, \rho'\}$.

First of all, for all $x \leq b_1^+$ and for all $t'' < t'$, $\mathcal{A}(x, b_2; \rho; t'') = \mathcal{A}(x, b_2; \rho'; t'')$, since otherwise the minimality of t' will be violated. The only difference in the allocation can occur in round t' .

Let us assume $\mathcal{A}_1(b_1, b_2; \rho; t') < \mathcal{A}_1(b_1, b_2; \rho'; t')$ (otherwise, we can swap ρ and ρ'). We make the claim that for all bids $x \leq b_1^+$ of agent 1, the influence of round t on round t' is in the same “direction”:

$$\mathcal{A}_1(x, b_2; \rho; t') \leq \mathcal{A}_1(x, b_2; \rho'; t') \text{ for all } x \leq b_1^+. \quad (3.2)$$

Suppose (3.2) does not hold. Then there is an $x < b_1^+$ such that $1 = \mathcal{A}_1(x, b_2; \rho; t') > \mathcal{A}_1(x, b_2; \rho'; t') = 0$. (Note that we have used the fact that the mechanism is deterministic.) If $x < b_1$ then pointwise monotonicity is violated under realization ρ , since $\mathcal{A}_1(x, b_2; \rho; t') > \mathcal{A}_1(b_1, b_2; \rho; t')$; otherwise it is violated under realization ρ' , giving a contradiction in both cases. The claim (3.2) follows.

Since \mathcal{A} is non-degenerate, there exists a non-degenerate interval I containing b_i such that if agent 1 bids any value $x \in I$ then $\mathcal{A}_1(x, b_2; \rho; t') < \mathcal{A}_1(x, b_2; \rho'; t')$. Now by (3.2) it follows that $I_1(\rho) < I_1(\rho')$. However, the mechanism cannot distinguish between ρ and ρ' when the bid of agent 1 is b_1^+ , since the differing bit $\rho_2(t)$ is not observed. Therefore the mechanism cannot compute prices, contradiction. \square

3.1 General Truthfulness Characterization

Let us develop the general truthfulness characterization that does not assume that an allocation is scale-free and IIA. We will later use it to derive Theorem 1.3.

Definition 3.4. Fix realization ρ and bid vector b . A round t is called $(b; \rho)$ -*secured* from agent i if $\mathcal{A}(b_i^+, b_{-i}; \rho; t) = \mathcal{A}(b_i, b_{-i}; \rho; t)$ for any $b_i^+ > b_i$. A round t is called *bid-independent* w.r.t. ρ if the allocation $\mathcal{A}(b; \rho; t)$ is a constant function of b .

The following definitions elaborate on the notion of an *influential round*.

Definition 3.5. A round t is called $(b; \rho)$ -*influential*, for bid profile b and realization ρ , if for some round $t' > t$ it holds that $\mathcal{A}(b; \rho; t') \neq \mathcal{A}(b; \rho'; t')$ for realization $\rho' = \rho \oplus \mathbf{1}(j, t)$ such that $j = \mathcal{A}(b; \rho; t)$.⁶ In this case, t' is called the *influenced round* and j is called the *influencing agent* of round t . The agent i is called an *influenced agent* of round t if $i \in \{\mathcal{A}(b; \rho; t'), \mathcal{A}(b; \rho'; t')\}$.

Note that a round is influential w.r.t. realization ρ if and only if it is (b, ρ) -influential for some b . The central property in our characterization is that each (b, ρ) -influential round is (b, ρ) -secured.

Definition 3.6. A deterministic allocation is called *weakly separated* if for every realization ρ and each bid vector b , it holds that if round t is $(b; \rho)$ -influential with influenced agent i then it is $(b; \rho)$ -secured from i .

We notice that exploration-separated is a stronger notion.

Observation 3.7. For a deterministic allocation, exploration-separated implies weakly separated.⁷

We are now ready to state our general characterization.

Theorem 3.8. Consider the MAB mechanism design problem. Let \mathcal{A} be a non-degenerate deterministic allocation rule. Then mechanism $(\mathcal{A}, \mathcal{P})$ is normalized and truthful for some payment rule \mathcal{P} if and only if \mathcal{A} is pointwise monotone and weakly separated.

Proof. For the “only if” direction, \mathcal{A} is pointwise-monotone by Lemma 3.2, and the fact that \mathcal{A} is weakly separated is proved similarly to Proposition 3.3 (albeit with a few extra details). We defer it to Appendix A.

We focus on the “if” direction. Let \mathcal{A} be a deterministic allocation rule which is pointwise monotone and weakly separated. We need to provide a payment rule \mathcal{P} such that the resulting mechanism $(\mathcal{A}, \mathcal{P})$ is truthful and normalized. Since \mathcal{A} is pointwise monotone, it immediately follows that it is monotone (i.e., as an agent increases her bid, the number of clicks that she gets cannot decrease). Therefore it follows from Theorem 3.1 that mechanism $(\mathcal{A}, \mathcal{P})$ is truthful and normalized if and only if \mathcal{P} is given by (3.1). We need to show that \mathcal{P} can be computed using only the knowledge of the clicks (bits from the realization) that were revealed during the execution of \mathcal{A} .

Assume we want to compute the payment for agent i in bid profile (b_i, b_{-i}) and realization ρ . We will prove that we can compute $\mathcal{C}_i(x) := \mathcal{C}_i(x, b_{-i}; \rho)$ for all $x \leq b_i$. To compute $\mathcal{C}_i(x)$, we show that it is possible to simulate the execution of the mechanism with $\text{bid}_i = x$. In some rounds, the agent i loses an impression, and in others it retains the impression (pointwise monotonicity ensures that agent i cannot gain an impression when decreasing her bid). In rounds that it loses an impression, the mechanism does not observe the bits of ρ in those rounds, so we prove that those bits are *irrelevant* while computing $\mathcal{C}_i(x)$. In other words, while running with $\text{bid}_i = x$, if mechanism needs to observe the bit that was not revealed when running with $\text{bid}_i = b_i$, we arbitrarily put that bit equal to 1 and simulate the execution of \mathcal{A} . We want to prove that this computes $\mathcal{C}_i(x)$ correctly.

Let $t_1 < t_2 < \dots < t_n$ be the rounds in which agent i did not get an impression while bidding x , but did get an impression while bidding b_i . Let $\rho^0 := \rho$, and let us define realization ρ^l inductively for every $l \in [n]$ by setting $\rho^l := \rho^{l-1} \oplus \mathbf{1}(j_l, t_l)$, where $j_l = \mathcal{A}(x, b_{-i}; \rho^{l-1}; t_l)$ is the agent that got the impression at round t_l with realization ρ^{l-1} and bids (x, b_{-i}) .

First, we claim that $j_l \neq i$ for any l . Indeed, suppose not, and pick the smallest l such that $j_{l+1} = i$. Then t_l is a $(x, b_{-i}; \rho^l)$ -influential round, with influenced agent $j_{l+1} = i$. Thus t_l is $(x, b_{-i}; \rho^l)$ -secured

⁶Note that realizations ρ and ρ' are interchangeable.

⁷To see this, simply use the definitions. Fix realization ρ and bid vector b , let t be a $(b; \rho)$ -influential round with influenced agent i . We need to show that t is $(b; \rho)$ -secured from i . Round t is $(b; \rho)$ -influential, thus influential w.r.t. ρ , thus (since the allocation is exploration-separated) it is bid-independent w.r.t. ρ , thus agent i cannot change allocation in round t by increasing her bid.

from i . Since $\mathcal{A}(x, b_{-i}; \rho^l; t_l) = \mathcal{A}(x, b_{-i}; \rho^{l-1}; t_l) = j_l \neq i$ by minimality of l , agent i does not get an impression in round t_l if she raises her bid to b_i . That is, $\mathcal{A}(b; \rho^l; t_l) \neq i$. However, the changes in realizations $\rho^0, \dots, \rho^{l-1}$ only concern the rounds in which agent i is chosen, so they are not seen by the allocation if the bid profile is b (to prove this formally, use induction). Thus, $\mathcal{A}(b; \rho^l; t_l) = \mathcal{A}(b; \rho; t_l) = i$, contradiction. Claim proved. It follows that $\mathcal{A}(b; \rho; t_l) = i$ for each l . (This is because by induction, the change from ρ^{l-1} to ρ^l is not seen by the allocation if the bid profile is b .)

We claim that $\mathcal{A}_i(x, b_{-i}; \rho; t') = \mathcal{A}_i(x, b_{-i}; \rho^n; t')$ for every round t' , which will prove the theorem. If not, then there exists l such that $\mathcal{A}_i(x, b_{-i}; \rho^l; t') \neq \mathcal{A}_i(x, b_{-i}; \rho^{l-1}; t')$ for some t' (and of course $t' > t_l$). Round t_l is thus $(x, b_{-i}; \rho^l)$ -influential with influenced round t' and influenced agent i . Moreover, the influencing agent of that round is j_l , and we already proved that $j_l \neq i$. Since round t_l is $(x, b_{-i}; \rho^l)$ -secured from agent i due to the “weakly separated” condition, it follows that agent i does not get an impression in round t_l if she raises her bid to b_i . That is, $\mathcal{A}(b; \rho^l; t_l) \neq i$, contradiction. \square

Note that we have proven the main characterization (Theorem 1.3) for the case of two agents, because for two agents IIA trivially holds and in the scale-free case, an allocation is exploration-separated if and only if it is weakly separated.

Let us argue that the non-degeneracy assumption in Theorem 3.8 is indeed necessary. To this end, let us present a simple deterministic mechanism $(\mathcal{A}, \mathcal{P})$ for two agents that is truthful and normalized, such that the allocation rule \mathcal{A} is pointwise monotone, scale-free and yet *not* weakly separated. (The catch is, of course, that it is degenerate.) There are only two rounds. Agent 1 allocated at round 1 if and only if $b_1 \geq b_2$. Agent 1 allocated at round 2 if $b_1 > b_2$ or if $b_1 = b_2$ and $\rho_1(1) = 1$; otherwise agent 2 is shown. This completes the description of the allocation rule. To obtain a payment rule \mathcal{P} which makes the mechanism normalized and truthful, consider an alternate allocation rule \mathcal{A}' which in each round selects agent 1 if and only if $b_1 \geq b_2$. (Note that $\mathcal{A}' = \mathcal{A}$ except when $b_1 = b_2$.) Use Theorem 3.8 for \mathcal{A}' to obtain a normalized truthful mechanism $(\mathcal{A}', \mathcal{P}')$, and set $\mathcal{P} = \mathcal{P}'$. The payment rule \mathcal{P} is well-defined since the observed clicks for \mathcal{P} and \mathcal{P}' coincide unless $b_1 = b_2$, in which case both payment rules charge 0 to both players. The resulting mechanism $(\mathcal{A}, \mathcal{P})$ is normalized and truthful because the integral in (3.1) remains the same even if we change the value at a single point. It is easy to see that the allocation rule \mathcal{A} has all the claimed properties; it fails to be non-degenerate because round t is influential only when $b_1 = b_2$.

3.2 Scalefree and IIA allocation rules

To complete the proof of Theorem 1.3, we show that under the right assumptions, an allocation is exploration-separated if and only if it is weakly separated. The full proof of this result is in Appendix A.

Lemma 3.9. *Consider the MAB mechanism design problem. Let \mathcal{A} be a non-degenerate deterministic allocation rule which is scalefree, pointwise monotone, and satisfies IIA. Then it is exploration-separated if and only if it is weakly separated.*

Proof Sketch. We sketch the proof of Lemma 3.9 at a very high level. The “only if” direction was observed in Observation 3.7. For the “if” direction, Let \mathcal{A} be a weakly-separated mechanism. We prove by a contradiction that it is exploration-separated. If not, then there is a realization ρ and a round t such that t is influential w.r.t. ρ as well as not bid-dependent w.r.t. ρ . Let round t be influential with bid vector b , influencing agent l , and influenced agents j and $j' \neq j$ in influenced round t' (see [1](#) in Figure 1; all boxed numbers in this sketch will refer to this figure).

From the assumption, t is not bid-dependent w.r.t. ρ , which means that there exists a bid profile b' such that $i' \neq l$ is played in round t with bids b' . Using scalefreeness, IIA, and pointwise-monotonicity, we can prove that there exists a sufficiently large bid $b_{i'}^+$ of agent i' such that she gets an impression in round t with

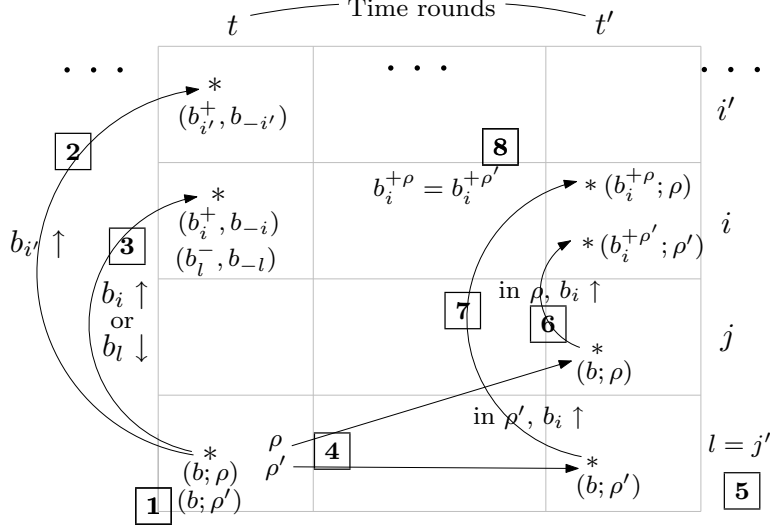


Figure 1: This figure explains all the steps in the proof of Lemma 3.9. The rows correspond to agents (whose identity is shown on the right side), and columns correspond to time rounds. The asterisks show the impressions. The arrows show how the impressions get *transferred*, and labels on the arrows show what causes the transfer. In labels, “in ρ , $b_i \uparrow$ ” denotes that a particular transfer of impression is caused in realization ρ when bid b_i is increased.

bids $(b_{i'}^+, b_{-i'})$ (see [2]). Using the properties of the mechanism, it can further be proved that there is an agent i such that she gets the impression in round t when either i increases her bid, *or* l decreases her bid (see [3]). When i increases her bid to b_i^+ , she also gets an impression in round t' , since impressions cannot differ in round t' in the case when l is not played in round t and they must get transferred from j and j' to *somebody* in round t' , and IIA implies that this *somebody* should be i .

Recall that two different players j and j' get the impression in round t' under ρ and ρ' respectively (see [4]). We prove that either agent j' or agent j must be equal to l (this is done by looking at how the allocation in round t' changes when l decreases her bid). Let us break the symmetry and assume $j' = l$ (see box [5]). It is also easy to see that when i increases her bid, impression in round t' get transferred to her in ρ (at some minimum value $b_i^{+\rho}$, see [6]), and impression in round t' gets transferred to her also in ρ' (as some possibly different minimum value $b_i^{+\rho'}$, see [7]). Using the assumptions of weakly-separatedness, we prove that $b_i^{+\rho} = b_i^{+\rho'}$ (see [8]). This can be proved by observing that $b_i^+ \geq \max\{b_i^{+\rho}, b_i^{+\rho'}\}$, and then using weakly-separatedness of \mathcal{A} . Since these two bids were at a “threshold value” (these were the minimum values of bids to have transferred the impression in ρ and ρ' from j and l respectively), we are able to prove that the ratio of b_j/b_l must be some fixed number dependent on ρ , ρ' , and t' . In particular, it follows that b_l belongs to a finite set $S(b_{-l})$ which depends only on b_{-l} . However, by non-degeneracy of \mathcal{A} there must be infinitely many such b_l ’s, which leads to a contradiction. \square

4 Lower bounds on regret

In this section we use structural results from the previous section to derive lower bounds on regret.

Theorem 4.1. *Consider the stochastic MAB mechanism design problem with k agents. Let \mathcal{A} be an exploration-separated deterministic allocation rule. Then its regret is $R(T; v_{\max}) = \Omega(v_{\max} k^{1/3} T^{2/3})$.*

Let $\vec{\mu}_0 = (\frac{1}{2}, \dots, \frac{1}{2}) \in [0, 1]^k$ be the vector of CTRs in which for each agent the CTR is $\frac{1}{2}$. For each agent i , let $\vec{\mu}_i = (\mu_{i1}, \dots, \mu_{ik}) \in [0, 1]^k$ be the vector of CTRs in which agent i has CTR $\mu_{ii} = \frac{1}{2} + \epsilon$,

$\epsilon = k^{1/3} T^{-1/3}$, and every other agent $j \neq i$ has CTR $\mu_{ij} = \frac{1}{2}$. As a notational convention, denote by $\mathbb{P}_i[\cdot]$ and $\mathbb{E}_i[\cdot]$ respectively the probability and expectation induced by the algorithm when clicks are given by $\vec{\mu}_i$. Let \mathcal{I}_i be the problem instance in which CTRs are given by $\vec{\mu}_i$ and all bids are v_{\max} . For each agent i , let \mathcal{J}_i be the problem instance in which CTRs are given by $\vec{\mu}_0$, the bid of agent i is v_{\max} , and the bids of all other agents are $v_{\max}/2$. We will show that for any exploration-separated deterministic allocation rule \mathcal{A} , one of these $2k$ instances causes high regret.

Let N_i be the number of bid-independent rounds in which agent i is played. Note that N_i does not depend on the bids. It is a random variable in the probability space induced by the clicks; its distribution is completely specified by the CTRs. We show that (in a certain sense) the allocation cannot distinguish between $\vec{\mu}_0$ and $\vec{\mu}_i$ if N_i is too small. Specifically, let \mathcal{A}_t be the allocation in round t . Once the bids are fixed, this is a random variable in the probability space induced by the clicks. For a given set S of agents, we consider the event $\{\mathcal{A}_t \in S\}$ for some fixed round t , and upper-bound the difference between the probability of this event under $\vec{\mu}_0$ and $\vec{\mu}_i$ in terms of $\mathbb{E}_i[N_i]$, in the following crucial claim, which is proved in Appendix B via relative entropy techniques.

Claim 4.2. *For any fixed vector of bids, each round t , each agent i and each set of agents S , we have*

$$|\mathbb{P}_0[\mathcal{A}_t \in S] - \mathbb{P}_i[\mathcal{A}_t \in S]| \leq O(\epsilon^2 \mathbb{E}_i[N_i]). \quad (4.1)$$

Proof of Theorem 4.1: Fix a positive constant β to be specified later. Consider the case $k = 2$ first. If $\mathbb{E}_0[N_i] > \beta T^{2/3}$ for some agent i , then on the problem instance \mathcal{J}_i , regret is $\Omega(T^{2/3})$. So without loss of generality let us assume $\mathbb{E}_0[N_i] \leq \beta T^{2/3}$ for each agent i . Then, plugging in the values for ϵ and $\mathbb{E}_0[N_i]$, the right-hand side of (4.1) is at most $O(\beta)$. Take β so that the right-hand side of (4.1) is at most $\frac{1}{4}$. For each round t there is an agent i such that $\mathbb{P}_0[\mathcal{A}_t \neq i] \geq \frac{1}{2}$. Then $\mathbb{P}_i[\mathcal{A}_t \neq i] \geq \frac{1}{4}$ by Claim 4.2, and therefore in this round algorithm \mathcal{A} incurs regret $\Omega(\epsilon v_{\max})$ under problem instance \mathcal{I}_i . By Pigeonhole Principle there exists an i such that this happens for at least half of the rounds t , which gives the desired lower-bound.

Case $k \geq 3$ requires a different (and somewhat more complicated) argument. Let $R = \beta k^{1/3} T^{2/3}$ and N be the number of bid-independent rounds. Assume $\mathbb{E}_0[N] > R$. Then $\mathbb{E}_0[N_i] \leq \frac{1}{k} \mathbb{E}_0[N]$ for some agent i . For the problem instance \mathcal{J}_i there are, in expectation, $E[N - N_i] = \Omega(R)$ bid-independent rounds in which agent i is not played; each of which contributes $\Omega(v_{\max})$ to regret, so the total regret is $\Omega(v_{\max} R)$.

From now on assume that $\mathbb{E}_0[N] \leq R$. Note that by Pigeonhole Principle, there are more than $\frac{k}{2}$ agents i such that $\mathbb{E}_0[N_i] \leq 2R/k$. Furthermore, let us say that an agent i is *good* if $\mathbb{P}_0[\mathcal{A}_t = i] \leq \frac{4}{5}$ for more than $T/6$ different rounds t . We claim that there are more than $\frac{k}{2}$ good agents. Suppose not. If agent i is not good then $\mathbb{P}_0[\mathcal{A}_t = i] > \frac{4}{5}$ for at least $\frac{5}{6}T$ different rounds t , so if there are at least $k/2$ such agents then

$$T = \sum_{t=1}^T \sum_{i=1}^k \mathbb{P}_0[\mathcal{A}_t = i] > \frac{k}{2} \times \left(\frac{5}{6}T\right) \times \frac{4}{5} \geq kT/3 \geq T,$$

contradiction. Claim proved. It follows that there exists a good agent i such that $\mathbb{E}_0[N_i] \leq 2R/k$. Therefore the right-hand side of (4.1) is at most $O(\beta)$. Pick β so that the right-hand side of (4.1) is at most $\frac{1}{10}$. Then by Claim 4.2 for at least $T/6$ different rounds t we have $\mathbb{P}_i[\mathcal{A}_t = i] \leq \frac{9}{10}$. In each such round, if agent i is not played then algorithm \mathcal{A} incurs regret $\Omega(\epsilon v_{\max})$ on problem instance \mathcal{I}_i . Therefore, the (total) regret of \mathcal{A} on problem instance \mathcal{I}_i is $\Omega(\epsilon v_{\max} T) = \Omega(v_{\max} k^{1/3} T^{2/3})$. \square

Theorem 4.3. *In the setting of Theorem 4.1, fix k and v_{\max} and assume that $R(T; v_{\max}) = O(v_{\max} T^\gamma)$ for some $\gamma < 1$. Then for every fixed $\delta \leq \frac{1}{4}$ and $\lambda < 2(1 - \gamma)$ we have $R_\delta(T; v_{\max}) = \Omega(\delta v_{\max} T^\lambda)$.*

Proof. Fix $\lambda \in (0, 2(1 - \gamma))$. Redefine $\vec{\mu}_i$'s with respect to a different ϵ , namely $\epsilon = T^{-\lambda/2}$. Define the problem instances \mathcal{I}_i in the same way as before: all bids are v_{\max} , the CTRs are given by $\vec{\mu}_i$.

Let us focus on agents 1 and 2. We claim that $\mathbb{E}_1[N_1] + \mathbb{E}_2[N_2] \geq \beta T^\lambda$, where $\beta > 0$ is a constant to be defined later. Suppose not. Fix all bids to be v_{\max} . For each round t , consider event $S_t = \{\mathcal{A}_t = 1\}$. Then by Claim 4.2 we have

$$|\mathbb{P}_1[S_t] - \mathbb{P}_2[S_t]| \leq |\mathbb{P}_0[S_t] - \mathbb{P}_1[S_t]| + |\mathbb{P}_0[S_t] - \mathbb{P}_2[S_t]| \leq O(\epsilon^2) (\mathbb{E}_1[N_1] + \mathbb{E}_2[N_2]) \leq \frac{1}{4}$$

for a sufficiently small β . Now, $\mathbb{P}_1[S_t] \geq \frac{1}{2}$ for at least $T/2$ rounds t . This is because otherwise on problem instance \mathcal{I}_1 regret would be $R(T) \geq \Omega(\epsilon T v_{\max}) = \Omega(v_{\max} T^{1-\lambda/2})$, which contradicts the assumption $R(T) = O(v_{\max} T^\gamma)$. Therefore $\mathbb{P}_2[S_t] \geq \frac{1}{4}$ for at least $T/2$ rounds t , hence on problem instance \mathcal{I}_2 regret is at least $\Omega(\epsilon T v_{\max})$, contradiction. Claim proved.

Now without loss of generality let us assume that $\mathbb{E}_1[N_1] \geq \frac{\beta}{2} T^\lambda$. Consider the problem instance in which CTRs given by $\vec{\mu}_1$, bid of agent 2 is v_{\max} , and all other bids are $v_{\max}(1 - 2\delta)/(1 + 2\epsilon)$. It is easy to see that this problem instance has δ -gap. Each time agent 1 is selected, algorithm incurs regret $\Omega(\delta v_{\max})$. Thus the total regret is at least $\Omega(\delta N_1 v_{\max}) = \Omega(\delta v_{\max} T^\lambda)$. \square

Matching upper bound. Let us describe a very simple mechanism, called *the naive MAB mechanism*, which matches the lower bound from Theorem 4.1 up to polylogarithmic factors (and also the lower bound from Theorem 4.3, for $\gamma = \lambda = \frac{2}{3}$ and constant δ).⁸

Fix the number of agents k , the time horizon T , and the bid vector b . The mechanism has two phases. In the *exploration phase*, each agent is played for $T_0 := k^{-2/3} T^{2/3} (\log T)^{1/3}$ rounds, in a round robin fashion. Let c_i be the number of clicks on agent i in the exploration phase. In the *exploitation phase*, an agent $i^* \in \arg\max_i c_i b_i$ is chosen and played in all remaining rounds. Payments are defined as follows: agent i^* pays $\max_{i \in [k] \setminus \{i^*\}} c_i b_i / c_{i^*}$ for every click she gets in exploitation phase, and all others pay 0. (Exploration rounds are free for every agent.) This completes the description of the mechanism.

Observation 4.4. *Consider the stochastic MAB mechanism design problem with k agents. The naive mechanism is normalized, truthful and has worst-case regret $R(T; v_{\max}) = O(v_{\max} k^{1/3} T^{2/3} \log^{2/3} T)$.*

Proof. The mechanism is truthful by a simple second-price argument.⁹ Recall that c_i is the number of clicks i got in the exploration phase. Let $p_i = \max_{j \neq i} c_j b_j / c_i$ be the price paid (per click) by agent i if she wins (all) rounds in exploitation phase. If $v_i \geq p_i$, then by bidding anything greater than p_i agent i gains $v_i - p_i$ utility each click irrespective of her bid, and bidding less than v_i , she gains 0, so bidding v_i is weakly dominant. Similarly, if $v_i < p_i$, then by bidding anything less than p_i she gains 0, while bidding $b_i > p_i$, she loses $b_i - p_i$ each click. So bidding v_i is weakly dominant in this case too.

For the regret bound, let (μ_1, \dots, μ_k) be the vector of CTRs, and let $\bar{\mu}_i = c_i / T_0$ be the sample CTRs. By Chernoff bounds, for each agent i we have $\Pr[|\bar{\mu}_i - \mu_i| > r] \leq T^{-4}$, for $r = \sqrt{8 \log(T) / T_0}$. If in a given run of the mechanism all estimates $\bar{\mu}_i$ lie in the intervals specified above, call the run *clean*. The expected regret from the runs that are not clean is at most $O(v_{\max})$, and can thus be ignored. From now on let us assume that the run is clean.

The regret in the exploration phase is at most $k T_0 v_{\max} = O(v_{\max} k^{1/3} T^{2/3} \log^{1/3} T)$. For the exploitation phase, let $j = \arg\max_i \mu_i b_i$. Then (since we assume that the run is clean) we have

$$(\mu_{i^*} + r) b_{i^*} \geq \bar{\mu}_{i^*} b_{i^*} \geq \bar{\mu}_j b_j \geq (\mu_j - r) b_j,$$

which implies $\mu_j v_j - \mu_{i^*} v_{i^*} \leq r(v_j + v_{i^*}) \leq 2r v_{\max}$. Therefore, the regret in exploitation phase is at most $2r v_{\max} T = O(v_{\max} k^{1/3} T^{2/3} \log^{2/3} T)$. Therefore the total regret is as claimed. \square

⁸Independently, Devanur and Kakade [17] presented a version of the naive MAB mechanism that achieves the same regret even in the more general model in which the value-per-click of an agent changes over time and the agents are allowed to submit a different bid at every round. Instead of assigning all impressions to the same agent in the exploitation phase, their mechanism runs the same allocation and payment procedure for each exploration round separately (see [17] for details).

⁹Alternatively, one can use Theorem 3.8 since all exploration rounds are bid-independent, and only exploration rounds are influential, and the payments are exactly as defined in Theorem 3.1.

5 Extensions and open questions

We extend our results in several directions. First, we derive a regret lower bound for deterministic truthful mechanisms without assuming that the allocations are scale-free. In particular, for two agents there are no assumptions. This lower bound holds for any k (the number of agents) assuming IIA, but unlike the one in Theorem 4.1 it does not depend on k . See Appendix C for details.

Second, we extend our results to randomized mechanisms. We consider randomized mechanisms that are *universally truthful*, i.e. truthful for each realization of the internal random seed. For mechanisms that randomize over exploration-separated deterministic allocation rules, we obtain the same lower bounds as in Theorems 4.1 and Theorem 4.3, see Appendix D for the details.

Third, we consider randomized allocation rules under a weaker version of truthfulness: a mechanism is *weakly truthful* if for each realization, it is truthful in expectation over its random seed. We show that any randomized allocation that is “pointwise monotone” and satisfies a certain notion of “separation between exploration and exploitation” can be turned into a mechanism that is weakly truthful and normalized. Then we apply this result to an algorithm in the literature [8, 25] in order to obtain regret guarantees for the version of the MAB mechanism design problem in which the clicks are chosen by an oblivious adversary.¹⁰ (This version corresponds to the *adversarial MAB problem* [7, 16, 1, 10].) The upper bound matches our lower bound for deterministic allocations up to $(\log k)^{1/3}$ factors. See Appendix E for details.

Fourth, we consider the stochastic MAB mechanism design problem under a more relaxed notion of truthfulness: truthfulness *in expectation*, where for each vector of CTRs the expectation is taken over clicks (and the internal randomness in the mechanism, if the latter is not deterministic). Following our line of investigation, we ask whether restricting a mechanism to be truthful in expectation has any implications on the structure and regret thereof. Given our results on mechanisms that are truthful and normalized, it is tempting to seek similar results for mechanisms that are truthful in expectation and normalized in expectation.¹¹ We rule out this approach: we show that in order to obtain any non-trivial lower bounds on regret and (essentially) any non-trivial structural results, one needs to assume that a mechanism is ex-post normalized, at least in some approximate sense. The key idea is to view the allocation and the payment as multivariate polynomials over the CTRs. See Appendix F for the details.

The two major open questions left open by this work concern structural results and regret lower bounds for (i) weakly truthful randomized mechanisms allocations, and (ii) mechanisms that are truthful in expectation. The latter question seems to require very different techniques which would further explore the connection to polynomials that we used in Appendix F.

Acknowledgements. We thank Jason Hartline, Bobby Kleinberg and Ilya Segal for helpful discussions.

References

- [1] Jacob Abernethy, Elad Hazan, and Alexander Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In *Conf. on Learning Theory (COLT)*, pages 263–274, 2008.
- [2] Gagan Aggarwal, Ashish Goel, and Rajeev Motwani. Truthful auctions for pricing search keywords. In *ACM Conf. on Electronic Commerce (EC)*, pages 1–7, 2006.
- [3] Gagan Aggarwal and S. Muthukrishnan. Tutorial on theory of sponsored search auctions. In *IEEE Symp. on Foundations of Computer Science (FOCS)*, 2008.

¹⁰An oblivious adversary chooses the entire realization in advance, without observing algorithm’s behavior and its random seed.

¹¹A mechanism is *normalized in expectation* if in expectation over clicks (and possibly over the allocation’s randomness), each agent is charged an amount between 0 and her bid for each click she receives.

- [4] Aaron Archer and Éva Tardos. Truthful mechanisms for one-parameter agents. In *IEEE Symp. on Foundations of Computer Science (FOCS)*, pages 482–491, 2001.
- [5] Susan Athey and Ilya Segal. An efficient dynamic mechanism. Available from <http://www.stanford.edu/~isegal/agv.pdf>, March 2007.
- [6] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2–3):235–256, 2002. Preliminary version in *15th ICML*, 1998.
- [7] Peter Auer, Nicolò Cesa-Bianchi, Yoav Freund, and Robert E. Schapire. The nonstochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002. Preliminary version in *36th IEEE FOCS*, 1995.
- [8] Baruch Awerbuch and Robert Kleinberg. Online linear optimization and adaptive routing. *Journal of Computer and System Sciences*, 74(1):97–114, February 2008. Preliminary version appeared in STOC 2004.
- [9] Moshe Babaioff, Yogeshwer Sharma, and Aleksandrs Slivkins. Characterizing truthful multi-armed bandit mechanisms. In *10th ACM Conf. on Electronic Commerce (EC)*, pages 79–88, 2009.
- [10] Peter L. Bartlett, Varsha Dani, Thomas Hayes, Sham Kakade, Alexander Rakhlin, and Ambuj Tewari. High-probability regret bounds for bandit online linear optimization. In *Conf. on Learning Theory (COLT)*, pages 335–342, 2008.
- [11] Michael Ben-Or and Avinatan Hassidim. The Bayesian Learner is Optimal for Noisy Binary Search (and Pretty Good for Quantum as Well). In *IEEE Symp. on Foundations of Computer Science (FOCS)*, 2008.
- [12] Dirk Bergemann and Juuso Välimäki. Efficient dynamic auctions. Available from cowles.econ.yale.edu/P/cd/d15b/d1584.pdf, October 2006.
- [13] Donald Berry and Bert Fristedt. *Bandit problems: sequential allocation of experiments*. Chapman&Hall, 1985.
- [14] Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- [15] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, New York, 1991.
- [16] Varsha Dani and Thomas P. Hayes. Robbing the bandit: Less regret in online geometric optimization against an adaptive adversary. In *16th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 937–943, 2006.
- [17] Nikhil Devanur and Sham M. Kakade. The price of truthfulness for pay-per-click auctions. In *10th ACM Conf. on Electronic Commerce (EC)*, pages 99–106, 2009.
- [18] Shahar Dobzinski and Mukund Sundararajan. On characterizations of truthful mechanisms for combinatorial auctions and scheduling. In *ACM Conf. on Electronic Commerce (EC)*, pages 38–47, 2008.
- [19] Benjamin Edelman, Michael Ostrovsky, and Michael Schwarz. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *American Economic Review*, 97(1):242–259, March 2007.
- [20] J. C. Gittins. *Multi-Armed Bandit Allocation Indices*. John Wiley & Sons, 1989.
- [21] Rica Gonen and Elan Pavlov. An incentive-compatible multi-armed bandit mechanism. In *ACM Symp. on Principles Of Distributed Computing (PODC) (Brief Announcement)*, pages 362–363, 2007. Preliminary version in *3rd Workshop on Sponsored Search Auctions* (in conjunction with WWW 2007).
- [22] Nicole Immorlica, Kamal Jain, Mohammad Mahdian, and Kunal Talwar. Click fraud resistant methods for learning click-through rates. In *Intl. Workshop On Internet And Network Economics (WINE)*, pages 34–45, 2005.
- [23] Richard Karp and Robert Kleinberg. Noisy binary search and its applications. In *18th ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 881–890, 2007.
- [24] Robert Kleinberg. *Online Decision Problems with Large Strategy Sets*. PhD thesis, MIT, Boston, MA, 2005.
- [25] Robert Kleinberg. Lecture notes: CS683: *Learning, Games, and Electronic Markets* (week 8), Spring 2007. Available at <http://www.cs.cornell.edu/courses/cs683/2007sp/lecnotes/week8.pdf>.

- [26] Robert Kleinberg. Lecture notes: *CS683: Learning, Games, and Electronic Markets* (week 9), Spring 2007. Available at <http://www.cs.cornell.edu/courses/cs683/2007sp/lecnotes/week9.pdf>.
- [27] Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. Multi-Armed Bandits in Metric Spaces. In *40th ACM Symp. on Theory of Computing (STOC)*, pages 681–690, 2008.
- [28] Sebastien Lahaie, David M. Pennock, Amin Saberi, and Rakesh V. Vohra. In *N. Nisan, T. Roughgarden, E. Tardos, and V. Vazirani (eds.) Chapter 28, Sponsored search auctions*. Cambridge University Press., 2007.
- [29] T.L. Lai and Herbert Robbins. Asymptotically efficient Adaptive Allocation Rules. *Advances in Applied Mathematics*, 6:4–22, 1985.
- [30] Ron Lavi, Ahuva Mu’alem, and Noam Nisan. Towards a characterization of truthful combinatorial auctions. In *IEEE Symp. on Foundations of Computer Science (FOCS)*, page 574, 2003.
- [31] Ron Lavi and Noam Nisan. Online ascending auctions for gradually expiring items. In *ACM-SIAM Symp. on Discrete Algorithms (SODA)*, pages 1146–1155, 2005.
- [32] Aranyak Mehta, Amin Saberi, Umesh Vazirani, and Vijay Vazirani. Adwords and generalized online matching. *J. ACM*, 54(5):22, 2007.
- [33] Roger B. Myerson. Optimal Auction Design. *Mathematics of Operations Research*, 6:58–73, 1981.
- [34] Hamid Nazerzadeh, Amin Saberi, and Rakesh Vohra. Dynamic cost-per-action mechanisms and applications to online advertising. In *17th Intl. World Wide Web Conf. (WWW)*, pages 179–188, 2008.
- [35] N. Nisan and A. Ronen. Algorithmic Mechanism Design. *Games and Economic Behavior*, 35(1-2):166–196, 2001.
- [36] N. Nisan, T. Roughgarden, E. Tardos, and V. Vazirani (eds.). *Algorithmic Game Theory*. Cambridge University Press., 2007.
- [37] Christos Papadimitriou, Michael Schapira, and Yaron Singer. On the hardness of being truthful. In *IEEE Symp. on Foundations of Computer Science (FOCS)*, 2008.
- [38] Tim Roughgarden. An algorithmic game theory primer. IFIP International Conference on Theoretical Computer Science (TCS). An invited survey., 2008.
- [39] Hal R. Varian. Position auctions. *International Journal of Industrial Organization*, 25(6):1163–1178, December 2007.

Appendix A: Truthfulness characterization

In this section we provide the proofs which did not fit into Section 3. We start with a complete proof of the “only if” direction of Theorem 3.8.

Lemma A.1. *Consider the MAB mechanism design problem. Let $(\mathcal{A}, \mathcal{P})$ be a normalized truthful mechanism such that \mathcal{A} is a non-degenerate deterministic allocation rule. Then \mathcal{A} is weakly separated.*

Proof. Assume \mathcal{A} is not weakly separated. Then there is a *counterexample* (ρ, b, t, t', i) : a realization ρ , bid vector b , rounds t, t' and agent i such that round t is $(b; \rho)$ -influential with influenced agent i and influenced round t' and it does not hold that round t is $(b; \rho)$ -secured from i . We prove that this leads to a contradiction..

Let us pick a counterexample (ρ, b, t, t', i) with a minimum value of t' over all choices of (ρ, b, t, i) . Without loss of generality, let us assume that $\rho_i(t') = 1$ and $\rho_j(t'') = 0$ for all $t'' > t'$ and for all agents j .

Let $j = \mathcal{A}(b; \rho; t)$. As it does not hold that round t is $(b; \rho)$ -secured from i , this means that $j \neq i$, and there exists a bid $b_i^+ > b_i$ such that $\mathcal{A}(b_i^+, b_{-i}; \rho; t) \neq j$.

Let $\rho' = \rho \oplus \mathbf{1}(j, t)$. The mechanism needs to compute prices for agent i when her bid is b_i^+ under realizations ρ and ρ' , that is to compute $\mathcal{P}_i(b_i^+, b_{-i}; \rho)$ and $\mathcal{P}_i(b_i^+, b_{-i}; \rho')$. Therefore, the mechanism needs to compute the integral $I_i(\varphi) = \int_{x \leq b_i^+} \mathcal{C}_i(x, b_{-i}; \varphi) dx$ for both realizations $\varphi \in \{\rho, \rho'\}$.

First of all, for all $x \leq b_i^+$ and for all $t'' < t'$, $\mathcal{A}_i(x, b_{-i}; \rho; t'') = \mathcal{A}_i(x, b_{-i}; \rho'; t'')$. If not, then the minimality of t' will be violated. This is because, if there were such an x and $t'' < t'$ with $\mathcal{A}_i(x, b_{-i}; \rho; t'') \neq \mathcal{A}_i(x, b_{-i}; \rho'; t'')$, then round t will still be (b, ρ) -influential with influenced agent i , and influenced round $t'' < t'$, violating the minimality of t' . Therefore, when we decrease the bid of agent i , the only difference in the allocation can occur at time round t' .

As i is the influenced agent at round t' it must hold that $\mathcal{A}_i(b_i, b_{-i}; \rho; t') \neq \mathcal{A}_i(b_i, b_{-i}; \rho'; t')$. Let us assume $0 = \mathcal{A}_i(b_i, b_{-i}; \rho; t') < \mathcal{A}_i(b_i, b_{-i}; \rho'; t') = 1$ (otherwise, we can swap ρ and ρ'). Note that we have made use of the fact that the mechanism is deterministic. Let us make the claim that for all bids $x \leq b_i^+$ the influence of round t on round t' is in the same “direction.”

$$\mathcal{A}_i(x, b_{-i}; \rho; t') \leq \mathcal{A}_i(x, b_{-i}; \rho'; t') \text{ for all } x \leq b_i^+. \quad (\text{A.1})$$

Suppose (A.1) does not hold. Then there is an $x \leq b_i^+$ such that $1 = \mathcal{A}_i(x, b_{-i}; \rho; t') > \mathcal{A}_i(x, b_{-i}; \rho'; t') = 0$. (Note that we have used the fact that the mechanism is deterministic.) If $x > b_i$, then pointwise monotonicity is violated in ρ' , since $0 = \mathcal{A}_i(x, b_{-i}; \rho'; t') < \mathcal{A}_i(b_i, b_{-i}; \rho'; t') = 1$. If $x < b_i$ on the other hand, then the pointwise-monotonicity is violated in ρ , since $1 = \mathcal{A}_i(x, b_{-i}; \rho; t') > \mathcal{A}_i(b_i, b_{-i}; \rho; t') = 0$, giving a contradiction in both cases. The claim (A.1) follows.

By the non-degeneracy of \mathcal{A} , there exists a non-degenerate interval I containing b_i such that

$$\mathcal{A}_i(x, b_{-i}; \rho; t') < \mathcal{A}_i(x, b_{-i}; \rho'; t') \text{ for all } x \in I. \quad (\text{A.2})$$

By (A.1) and (A.2) it follows that $I_i(\rho) < I_i(\rho')$. However, the mechanism cannot distinguish between ρ and ρ' when agent i 's bid is b_i^+ , since the differing bit $\rho_j(t)$ is not seen. Contradiction. \square

A.1 Proof of Lemma 3.9

For convenience, let us restate the lemma.

Lemma (Lemma 3.9 restated). *Consider the MAB mechanism design problem. Let \mathcal{A} be a non-degenerate deterministic allocation rule which is scalefree, pointwise monotone, and satisfies IIA. Then it is exploration-separated if and only if it is weakly separated.*

The “only if” direction is a consequence of Observation 3.7. Here we prove the “if” direction. For bid profile b , realization ρ , agent l and round t , the tuple $(b; \rho; l; t)$ is called an *influence-tuple* if round t is (b, ρ) -influential with influencing agent l . Suppose allocation \mathcal{A} is weakly separated but not exploration-separated. Then there is a *counterexample*: an influence-tuple $(b; \rho; l; t)$ such that round t is not bid-independent w.r.t. realization ρ . We prove that such counterexample can occur only if $b_l \in S_l(b_{-l})$, for some finite set $S_l(b_{-l}) \subset \mathbb{R}$ that depends only on b_{-l} .

Proposition A.2. *Let \mathcal{A} be as in Lemma 3.9. Assume \mathcal{A} is weakly separated. Then for each agent l and each bid profile b_{-l} there exists a finite set $S_l(b_{-l}) \subset \mathbb{R}$ with the following property: for each counterexample $(b_l, b_{-l}; \rho; l; t)$ it is the case that $b_l \in S_l(b_{-l})$.*

Once this proposition is proved, we obtain a contradiction with the non-degeneracy of \mathcal{A} . Indeed, suppose $(b; \rho; l; t)$ is a counterexample. Then $(b; \rho; l; t)$ is an influence-tuple. Since \mathcal{A} is non-degenerate, there exists a non-degenerate interval I such that for each $x \in I$ it holds that $(x, b_{-l}; \rho; l; t)$ is an influence-tuple, and therefore a counterexample. Thus the set $S_l(b_{-l})$ in Proposition A.2 cannot be finite, contradiction.

In the rest of this section we prove Proposition A.2. Fix a counterexample $(b; \rho; l; t)$; let $t' > t$ be the influenced round. In particular, $\mathcal{A}(b; \rho; t) = l$ (see [1](#) in Figure 1 on page 11; all boxed numbers will refer to this figure). Then by the assumption there exist bids b' such that $\mathcal{A}(b'; \rho; t) = i' \neq l$. We claim

that this implies that there exists a bid $b_{i'}^+ > b_{i'}$ such that $\mathcal{A}(b_{i'}^+, b_{-i'}; \rho; t) = i'$ (see [2]). This is proven in Lemma A.4 below, and in order to prove it we first present the following lemma, which essentially states that if the mechanism makes a choice between i and j of who to be show, then it can only depend on the ratio of their bids $\text{bid}_i/\text{bid}_j$, and not on the bids of other agents.

Lemma A.3. *Let \mathcal{A} be an MAB (deterministic) allocation rule that is pointwise-monotone, scalefree, and satisfies IIA. Let there be two bid profiles α and β such that $\mathcal{A}(\alpha; \rho; t) \in \{i, j\}$, $\mathcal{A}(\beta; \rho; t) \in \{i, j\}$, and $\alpha_i/\alpha_j = \beta_i/\beta_j$. Then it must be the case that $\mathcal{A}(\alpha; \rho; t) = \mathcal{A}(\beta; \rho; t)$.*

Proof. As \mathcal{A} is scalefree we assume that $\alpha_i = \beta_i$ and $\alpha_j = \beta_j$ by scaling bids in β by a factor of α_i/β_i (or a factor of α_j/β_j), without changing the allocation.

Assume for the sake of a contradiction that $\mathcal{A}(\beta; \rho; t) \neq \mathcal{A}(\alpha; \rho; t)$. Let us number the agents as follows. Agents i and j are numbered 1 and 2, respectively. The rest of the agents are arbitrarily numbered 3 to k . Consider the following sequence of bid vectors. $\alpha(1) = \alpha(2) = \alpha$ and $\alpha(m) = (\beta_m, \alpha(m-1)_{-m})$ for $m \in \{3, \dots, k\}$. As $\alpha(1) = \alpha$ and $\alpha(k) = \beta$, $\mathcal{A}(\alpha(1); \rho; t) = \mathcal{A}(\alpha; \rho; t)$ and $\mathcal{A}(\alpha(k); \rho; t) = \mathcal{A}(\beta; \rho; t)$. Since $\mathcal{A}(\alpha(k); \rho; t) = \mathcal{A}(\beta; \rho; t) \neq \mathcal{A}(\alpha; \rho; t) = \mathcal{A}(\alpha(1); \rho; t)$ there exists $m \in \{3, \dots, k\}$ such that $\mathcal{A}(\alpha(m-1); \rho; t) = \mathcal{A}(\alpha; \rho; t) \in \{i, j\}$ while $\mathcal{A}(\alpha(m); \rho; t) \neq \mathcal{A}(\alpha(m-1); \rho; t)$. As $m \neq i$ and $m \neq j$, IIA implies that $\mathcal{A}(\alpha(m); \rho; t) = m$ and given that, IIA also implies that $\mathcal{A}(\alpha(k); \rho; t) \in \{m, m+1, \dots, k\}$ (note that i, j are not in this set). But as $\mathcal{A}(\alpha(k); \rho; t) = \mathcal{A}(\beta; \rho; t) \in \{i, j\}$ this yields a contradiction. \square

Lemma A.4. *Let \mathcal{A} be an MAB (deterministic) allocation rule that is pointwise-monotone, scalefree, and satisfies IIA. Let there be two bid profiles α and β such that $\mathcal{A}(\alpha; \rho; t) = i$ and $\mathcal{A}(\beta; \rho; t) = j \neq i$. Then there exists $\beta_i^+ > \beta_i$ such that $\mathcal{A}(\beta_i^+, \beta_{-i}; \rho; t) = i$.*

In other words, if it is possible for i to get the impression in round t at all, then it is possible for her to get the impression starting from any bid profile and raising her bid high enough.

Proof. We first note that $\frac{\alpha_i}{\alpha_j} \geq \frac{\beta_i}{\beta_j}$. If not, then $\frac{\alpha_i}{\alpha_j} < \frac{\beta_i}{\beta_j}$. Consider a raised bid of i from α_i to $\alpha_i^+ = \alpha_j \cdot \frac{\beta_i}{\beta_j}$. In the bid profile $(\alpha_i^+, \alpha_{-i})$, i must get the impression (by pointwise monotonicity). This gives a contradiction to Lemma A.3, since $\mathcal{A}(\alpha_i^+, \alpha_{-i}; \rho; t) = i \in \{i, j\}$, $\mathcal{A}(\beta; \rho; t) = j \in \{i, j\}$, and $\frac{\alpha_i^+}{\alpha_j} = \frac{\beta_i}{\beta_j}$, but $\mathcal{A}(\alpha_i^+, \alpha_{-i}; \rho; t) \neq \mathcal{A}(\beta; \rho; t)$.

Now, consider i increasing her bid in profile β to $\beta_i^+ = \beta_j \cdot \frac{\alpha_i}{\alpha_j}$. Now, $\mathcal{A}(\alpha; \rho; t) = i \in \{i, j\}$, $\mathcal{A}(\beta_i^+, \beta_{-i}; \rho; t) \in \{i, j\}$ (from IIA), and $\frac{\alpha_i}{\alpha_j} = \frac{\beta_i^+}{\beta_j}$. We can apply Lemma A.3 to deduce that $\mathcal{A}(\alpha; \rho; t) = \mathcal{A}(\beta_i^+, \beta_{-i}; \rho; t)$ and both are equal to i since the first allocation is equal to i . \square

From the lemma above, it follows that agent i' can increase her bid (in bid profile b) and get the impression in realization ρ , round t . To quantify by how much agent i' needs to raise her bid to get the impression, we introduce the notion of *threshold* $\Theta_{i,j}(\rho; t)$ in the next lemma.

Lemma A.5. *Let \mathcal{A} be an MAB (deterministic) allocation rule that is pointwise monotone, scalefree and satisfies IIA. For realization ρ , round t , two agents i and $j \neq i$, let bids b_{-i-j} be such that there exist x_0 and y satisfying $\mathcal{A}(x_0, y, b_{-i-j}; \rho; t) = j$, and there exists x (possibly dependent on y) satisfying $\mathcal{A}(x, y, b_{-i-j}; \rho; t) = i$. Let us fix such a y and define¹²*

$$\Theta_{i,j}^{b_{-i-j}}(\rho, t) = \frac{1}{y} \inf_x \{x \mid \mathcal{A}(x, y, b_{-i-j}; \rho; t) = i\}.$$

¹²Note that if there are no values of bids of i (x_0 and x) and j (equal to y) such that j can get an impression with small enough bid (x_0) of agent i and i can get an impression by raising her bid (to x), then we don't define $\Theta_{i,j}^{b_{-i-j}}(\rho; t)$ at all. We will be careful not to use such undefined Θ 's. It is not hard to see that if bids are nonzero, then $\Theta_{i,j}(\rho; t)$ is defined if and only if $\Theta_{j,i}(\rho; t)$ is. Moreover $0 < \Theta_{i,j}(\rho; t) < \infty$, and $\Theta_{j,i}(\rho; t) = (\Theta_{i,j}(\rho; t))^{-1}$.

Then for any bids b'_{-i-j} , $\Theta_{i,j}^{b'-i-j}(\rho, t)$ is well defined and satisfies $\Theta_{i,j}^{b'-i-j}(\rho, t) = \Theta_{i,j}^{b-i-j}(\rho, t)$. We denote it by $\Theta_{i,j}(\rho, t)$, as $\Theta_{i,j}^{b-i-j}(\rho, t)$ is independent of b_{-i-j} .

Proof. We first prove that if the conditions of the definition of $\Theta_{i,j}^{b-i-j}(\rho; t)$ are satisfied for b_{-i-j} , then are also satisfied for any other b'_{-i-j} . Let us say they are satisfied for b_{-i-j} , that is there exists x_0, x and y , such that $\mathcal{A}(x_0, y, b_{-i-j}; \rho; t) = j$ and $\mathcal{A}(x, y, b_{-i-j}; \rho; t) = i$. We want to prove existence of x' and y' for b'_{-i-j} . If $\mathcal{A}(x_0, y, b'_{-i-j}; \rho; t) = j$ then existence of y' is proved for b'_{-i-j} too, since $y' = y$ works. If not, then $\mathcal{A}(x_0, y, b'_{-i-j}; \rho; t) = j' \neq j$ and $\mathcal{A}(x_0, y, b_{-i-j}; \rho; t) = j$, and by Lemma A.4, there exists a $y' > y$ such that $\mathcal{A}(x_0, y', b'_{-i-j}; \rho; t) = j$. Once the existence of y' is proved, we now prove the existence of x' . Let $x' = x \cdot \frac{y'}{y} \geq x$. We have $\mathcal{A}(x, y, b_{-i-j}; \rho; t) = i \in \{i, j\}$ and $\mathcal{A}(x', y', b'_{-i-j}; \rho; t) \in \{i, j\}$ by IIA (i can only transfer impression to her by changing her bid) and $x'/y' = x/y$. From Lemma A.3, we get $i = \mathcal{A}(x, y, b_{-i-j}; \rho; t) = \mathcal{A}(x', y', b'_{-i-j}; \rho; t)$. Hence the existence of x' is proved too.

For the sake of contradiction, let us assume that $\theta := \Theta_{i,j}^{b-i-j}(\rho; t) < \Theta_{i,j}^{b'-i-j}(\rho; t) =: \theta'$. Let us scale the bids in (x', y', b'_{-i-j}) by a factor such that the factor times y' is equal to y . We can hence assume that $y' = y$. Let us pick a bid $x'' \in (\theta y, \theta' y)$. We have $\mathcal{A}(x'', y, b_{-i-j}; \rho; t) = i$ (since x''/y is past the threshold θ), $\mathcal{A}(x'', y' = y, b'_{-i-j}; \rho; t) = j$ (x''/y' is yet not past the threshold θ'), and $x''/y = x''/y'$. This is a contradiction to the Lemma A.3. Therefore, $\theta = \theta'$. \square

We conclude that if $b_{i'}^+ > b_l \cdot \Theta_{i',l}(\rho, t)$ then $\mathcal{A}(b_{i'}^+, b_{-i'}; \rho; t) = i' \neq l$ (see [2] again). Note that we are using $\Theta_{i',l}(\rho; t)$ since this is well-defined. Define $\rho' = \rho \oplus \mathbf{1}(l, t)$.

Let us think about decreasing the bid of agent l from b_l (it is positive, since all bids are assumed to be positive). When the bid of agent l is b_l , she gets the impression in round t , but when her bid is small enough (in particular as low as $b_{i'}/\Theta_{i',l}(\rho; t)$), then she must not get the impression in round t (see Lemma A.3). When the bid of l decreases, some other agent gets the impression in round t , let us call that agent i (note that this agent may not be the same as agent i' above). See [3].

Now, starting from bid profile b , let us increase the bid of agent i . When the bid of agent i is large enough (in particular as large as $b_i \Theta_{i',l}(\rho; t) b_l / b_{i'}$), then l can no longer get the impression in round t (see Lemma A.3). From IIA, the impression must get transferred to i . Therefore we can define $\Theta_{i,l}(\rho; t)$, and when $b_i^+ > b_l \Theta_{i,l}(\rho; t)$, agent i gets the impression in round t (see [3] again). Note that $\mathcal{A}(b_i^+, b_{-i}; \rho; t) = \mathcal{A}(b_i^+, b_{-i}; \rho'; t) = i$ (click information for l at round t cannot influence the impression decision at round t).

Recall that t' is the influenced round. Let $\mathcal{A}(b; \rho; t') = j$ and let $\mathcal{A}(b; \rho'; t') = j' \neq j$ (see [4]). As \mathcal{A} is pointwise monotone and IIA, $\mathcal{A}(b_i^+, b_{-i}; \rho; t') \in \{i, j\}$ and $\mathcal{A}(b_i^+, b_{-i}; \rho'; t') \in \{i, j'\}$. It must be the case that $\mathcal{A}(b_i^+, b_{-i}; \rho; t') = \mathcal{A}(b_i^+, b_{-i}; \rho'; t')$, as l does not get an impression at round t (and the algorithm does not see the difference between ρ and ρ'). As $j' \neq j$ we conclude that

$$\mathcal{A}(b_i^+, b_{-i}; \rho; t') = \mathcal{A}(b_i^+, b_{-i}; \rho'; t') = i.$$

Next we note that $i \neq j$ and $i \neq j'$. This is because if $i = j$ (respectively $i = j'$), then round t would be $(b; \rho)$ -influential (respectively $(b; \rho')$ -influential) with influenced agent i but it is not $(b; \rho)$ -secured (respectively $(b; \rho')$ -secured) from i , in contradiction to the assumption.

We also note that $l \in \{j, j'\}$ (see [5]). Assume for the sake of contradiction that $l \neq j$ and $l \neq j'$. For $b_l^- < b_i \cdot \Theta_{l,i}(\rho, t)$ it holds that $\mathcal{A}(b_l^-, b_{-l}; \rho; t) = \mathcal{A}(b_l^-, b_{-l}; \rho'; t) = i$ (since i was defined such that i gets the impression in round t when l decreases her bid) thus $\mathcal{A}(b_l^-, b_{-l}; \rho; t') = \mathcal{A}(b_l^-, b_{-l}; \rho'; t')$ (as click information for l at round t is not observed). (Also, as a side note, observe that $b_l^- < b_l$ by pointwise-monotonicity since agent l was getting an impression in round t with bid b_l and lost it when her bid is b_l^- .) Let $\mathcal{A}(b_l^-, b_{-l}; \rho; t') = \mathcal{A}(b_l^-, b_{-l}; \rho'; t') = l'$. Note that $l' \neq l$, since otherwise, $\mathcal{A}_l(x, b_{-l}; \rho; t')$ is not a monotone function of x : it is 0 when $x = b_l$ (since j gets an impression), and 1 when $x = b_l^- < b_l$, a contradiction to pointwise-monotonicity. Now, note that the impression in ρ' at time t' transfers from j' to

l' , and impression in ρ at time t' transfers from j to l' , none of which $(\{j, j', l'\})$ are equal to l and $j \neq j'$. Let us write this in equations:

$$\begin{aligned} \mathcal{A}(b_l, b_{-l}; \rho; t') &= j & \mathcal{A}(b_l^-, b_{-l}; \rho; t') &= l' \\ \mathcal{A}(b_l, b_{-l}; \rho'; t') &= j' & \mathcal{A}(b_l^-, b_{-l}; \rho'; t') &= l'. \end{aligned}$$

It must be the case that either $j \neq l'$ or $j' \neq l'$ (since $j \neq j'$). If $j \neq l'$, then in ρ at time t' , reducing the bid of l transfers impression from j to l' (both of them are different from l), thus violating IIA. Similarly, if $j' \neq l'$, then in ρ' at time t' , reducing the bid of l transfers impression from j' to l' (both of them are different from l), thus violating IIA. We thus have $l \in \{j, j'\}$. Let $l = j'$ (since otherwise, we can swap the roles of ρ and ρ').

To summarize what we have proved so far: there are 3 distinct agents i, j, l such that

$$\begin{aligned} \mathcal{A}(b; \rho; t) &= \mathcal{A}(b; \rho'; t) = \mathcal{A}(b; \rho'; t') = l \quad (\text{since } \mathcal{A}(b; \rho'; t') = j' = l), \\ \mathcal{A}(b; \rho; t') &= j \quad \text{and} \\ \mathcal{A}(b_i^+, b_{-i}; \rho; t) &= \mathcal{A}(b_i^+, b_{-i}; \rho; t') = \mathcal{A}(b_i^+, b_{-i}; \rho'; t) = \mathcal{A}(b_i^+, b_{-i}; \rho'; t') = i. \end{aligned}$$

Observe also that $\Theta_{i,l}(\rho, t) = \Theta_{i,l}(\rho', t)$ as ρ and ρ' only differ at a click at round t , and such a click cannot determine the allocation decision at round t . Also, $\max\{\Theta_{i,j}(\rho, t') \cdot b_j, \Theta_{i,l}(\rho', t') \cdot b_l\} \leq \Theta_{i,l}(\rho, t) \cdot b_l$ as the allocation at round t' , which is different for ρ and ρ' (at b), depends on l getting the impression at round t .¹³ Finally we prove that $\Theta_{i,j}(\rho, t') \cdot b_j = \Theta_{i,l}(\rho', t') \cdot b_l$ (see [8]).

Claim A.6. $\Theta_{i,j}(\rho, t') \cdot b_j = \Theta_{i,l}(\rho', t') \cdot b_l$

Proof. First of all, note that $\Theta_{i,j}(\rho, t')$ and $\Theta_{i,l}(\rho', t')$ are well-defined. Let $\bar{b}_i = (\Theta_{i,j}(\rho, t') \cdot b_j + \Theta_{i,l}(\rho', t') \cdot b_l) / 2$. Consider the following two cases.

If $\Theta_{i,j}(\rho, t') \cdot b_j < \Theta_{i,l}(\rho', t') \cdot b_l$ then round t is $(\bar{b}_i, b_{-i}; \rho)$ -influential (as $\mathcal{A}(\bar{b}_i, b_{-i}; \rho; t') = i$ and $\mathcal{A}(\bar{b}_i, b_{-i}; \rho'; t') = l$) with influencing agent l ($\mathcal{A}(\bar{b}_i, b_{-i}; \rho; t) = \mathcal{A}(\bar{b}_i, b_{-i}; \rho'; t) = l$ since $\bar{b}_i < \Theta_{i,l}(\rho, t) \cdot b_l$) and influenced agent i . Additionally, t is not $(\bar{b}_i, b_{-i}; \rho)$ -secured from i (as $\mathcal{A}(b_i^+, b_{-i}; \rho; t) = \mathcal{A}(b_i^+, b_{-i}; \rho'; t) = i$). A contradiction to first condition in the theorem.

Similarly, if $\Theta_{i,j}(\rho, t') \cdot b_j > \Theta_{i,l}(\rho', t') \cdot b_l$ then round t is $(\bar{b}_i, b_{-i}; \rho)$ -influential (as now $\mathcal{A}(\bar{b}_i, b_{-i}; \rho; t') = j$ and $\mathcal{A}(\bar{b}_i, b_{-i}; \rho'; t') = i$) with influencing agent l and influenced agent i . Additionally, t is not $(\bar{b}_i, b_{-i}; \rho)$ -secured from i . Again, a contradiction to the first condition in the theorem. \square

The lemma implies that $b_l \in S_l(b_{-l})$, where a finite set $S_l(b_{-l})$ is defined by

$$S_l(b_{-l}) = \left\{ b_j \frac{\Theta_{i,j}(\rho, t')}{\Theta_{i,l}(\rho', t')} : \text{all agents } i, j \neq l, \text{ all realizations } \rho, \rho' \text{ and all } t' \text{ s.t. } \frac{\Theta_{i,j}(\rho, t')}{\Theta_{i,l}(\rho', t')} \text{ is well-defined} \right\}.$$

This completes the proof of Proposition A.2.

Appendix B: Relative entropy technique: proof of Claim 4.2

We extend the relative entropy technique from [7]. All relevant facts about relative entropy are summarized in the theorem below. We will need the following definition: given a random variable X on a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, let \mathbb{P}_X be the distribution of X , i.e. a measure on \mathbb{R} defined by $\mathbb{P}_X(x) = \mathbb{P}[X = x]$.

Theorem B.1. *Let p and q be two probability measures on a finite set U , and let Y and Z be functions on U . There exists a function $F(p; q|Y) : U \rightarrow \mathbb{R}$ with the following properties:*

¹³In Figure 1 we defined $b_i^{+\rho} := \Theta_{i,j}(\rho; t')b_j$ and $b_i^{+\rho'} := \Theta_{i,l}(\rho'; t')b_l$. These are the bids of agent i at which impression transfers to her in round t' in ρ and ρ' respectively. See [6] and [7] in the figure.

- (i) $E_p F(p; q|Y) = E_p F(p; q|(Y, Z)) + E_p F(p_Z; q_Z|Y)$ (chain rule),
- (ii) $|p(U') - q(U')| \leq \sqrt{\frac{1}{2}\mathcal{D}(p||q)}$ for any event $U' \subset U$, where $\mathcal{D}(p||q) = E_p F(p; q|1)$
- (iii) for each $x \in U$, if conditional on the event $\{Z = Z(x)\}$ p coincides with q , then $F(p; q|Z)(x) = 0$.
- (iv) for each $x \in U$, if conditional on the event $\{Z = Z(x)\}$ p and q are fair and $(\frac{1}{2} + \epsilon)$ -biased coins, respectively, then it is the case that $F(p; q|Z)(x) \leq 4\epsilon^2$.

Remark. This theorem summarizes several well-known facts about relative entropy (albeit in a somewhat non-standard notation). For the proofs, see [15, 24, 26]. In the proofs, one defines $F = F(p; q|Y)$ as a function $F : U \rightarrow \mathbb{R}$ which is specified by $F(x) = \sum_{x' \in U} p(x'|U_x) \lg \frac{p(x'|U_x)}{q(x'|U_x)}$, where U_x is the event $\{Y = Y(x)\}$.¹⁴ Note that the quantity $E_p F(p; q|1)$ is precisely the relative entropy (a.k.a. KL-divergence), commonly denoted $\mathcal{D}(p||q)$, and $E_p F(p; q|Y)$ is the corresponding conditional relative entropy.

In what follows we use Theorem B.1 to prove Claim 4.2. For simplicity we will prove (4.1) for $i = 1$.

The *history* up to round t is $H_t = (h_1, h_2, \dots, h_t)$ where $h_s \in \{0, 1\}$ is the click or no click event received by the algorithm at round s . Let C_t be the indicator function of the event “round t is bid-independent”. Define the *bid-independent history* as $\hat{H}_t = (\hat{h}_1, \hat{h}_2, \dots, \hat{h}_t)$, where $\hat{h}_t = h_t C_t$. For any exploration-separated deterministic allocation rule and each round t , the bid-independent history \hat{H}_{t-1} and the bids completely determine which arm is chosen in this round. Moreover, \hat{H}_{t-1} alone (without the bids) completely determines whether round t is bid-independent, and if so, which arm is chosen in this round.

Recall the CTR vectors $\vec{\mu}_i$ as defined in Section 4. Let p and q be the distributions induced on \hat{H}_T by $\vec{\mu}_0$ and $\vec{\mu}_1$, respectively. Let p_t and q_t be the distributions induced on \hat{h}_t by $\vec{\mu}_0$ and $\vec{\mu}_1$, respectively. Let \mathcal{H}_t the support of \hat{H}_t , i.e. the set of all t -bit vectors. In the forthcoming applications of Theorem B.1, the universe will be $U = \mathcal{H}_T$. By abuse of notation, we will treat \hat{H}_t as a projection $\mathcal{H}_T \rightarrow \mathcal{H}_t$, so that it can be considered a random variable under p or q .

Claim B.2. $\mathcal{D}(p||q) = E_p F(p; q|\hat{H}_t) + \sum_{s=1}^t E_p F(p_s; q_s|\hat{H}_{s-1})$ for any $t > 1$.

Proof. Use induction on $t \geq 0$ (set $\hat{H}_0 = 1$). In order to obtain the claim for a given t assuming that it holds for $t - 1$, apply Theorem B.1(i) with $Y = \hat{H}_{t-1}$ and $Z = \hat{h}_t$. \square

Claim B.3. $F(p_t; q_t|\hat{H}_{t-1}) \leq 4\epsilon^2 C_t 1_{\{A_t=1\}}$ for each round t .

Proof. We are interested in the function $F = F(p_t; q_t|\hat{H}_{t-1}) : \mathcal{H}_T \rightarrow \mathbb{R}$. Given \hat{H}_{t-1} , one of the following three cases occurs:

- round t is not bid-independent. Then $\hat{h}_t = 0$, hence $F(\cdot) = 0$ by Theorem B.1(iii),
- round t is bid-independent and arm 1 is not played. Then \hat{h}_t is distributed as a fair coin under both p and q , so again $F(\cdot) = 0$.
- round t is bid-independent and arm 1 is played. Then $F(\cdot) \leq 4\epsilon^2$ by Theorem B.1(iv). \square

Given the full bid-independent history \hat{H}_T , p and q become (the same) point measure, so by Theorem B.1(iii) $E_p F(p; q|\hat{H}_T) = 0$. Therefore taking Claim B.2 with $t = T$ we obtain

$$\mathcal{D}(p||q) = \sum_{t=1}^T E_p F(p_t; q_t|\hat{H}_{t-1}) = 4\epsilon^2 \sum_{t=1}^T E_p [C_t 1_{\{A_t=1\}}] = 4\epsilon^2 E_p [N_1]. \quad (\text{B.1})$$

For a given round t and fixed bids, the allocation at round t is completely determined by the bid-independent history \hat{H}_{t-1} . Thus, we can treat $\{A_t \in S\}$ as an event in \mathcal{H}_T . Now (4.1) follows from (B.1) via an application of Theorem B.1(ii) with $U' = \{A_t \in S\}$.

¹⁴We use the convention that $p(x) \log(p(x)/q(x))$ is 0 when $p(x) = 0$, and $+\infty$ when $p(x) > 0$ and $q(x) = 0$.

Appendix C: Lower bound for non-scalefree allocations

In this section we derive a regret lower bound for deterministic truthful mechanisms without assuming that the allocations are scale-free. In particular, for two agents there are no assumptions. This lower bound holds for any k (the number of agents) assuming that the allocation satisfies IIA, but unlike the one in Theorem 4.1 it does not depend on k .

Theorem C.1. *Consider the stochastic MAB mechanism design problem with k agents. Let $(\mathcal{A}, \mathcal{P})$ be a normalized truthful mechanism such that \mathcal{A} is a non-degenerate deterministic allocation rule. Suppose \mathcal{A} satisfies IIA. Then its regret is $R(T; v_{\max}) = \Omega(v_{\max} T^{2/3})$ for any sufficiently large v_{\max} .*

Let us sketch the proof. Fix an allocation \mathcal{A} . In Definition 3.4, if round t is (b, ρ) influential, for some realization ρ and bid vector b , an agent i is called *strongly influenced* by round t if it is one of the two agents that are “influenced” by round t but is not the “influencing agent” of round t . In particular, it holds that $\mathcal{A}(b, \rho, t) \neq i$. For each realization ρ , round t and agent i , if there exists a bid vector b such that round t is (b, ρ) -influential with strongly influenced agent i , then fix any one such b , and define $b_i^* = b_i^*(\rho, t) := \max_{j \neq i} b_j$. Let us define $B_{\mathcal{A}}^* = \max_{\rho, t, i} b_i^*(\rho, t)$, where the maximum is taken over all realizations ρ , all rounds t , and all agents i . Let us say that round t is B^* -free from agent i w.r.t realization ρ , if for this realization the following property holds: agent i is not selected in round t as long as each bid is at least B^* .

Lemma C.2. *In the setting of Theorem C.1, for any realization ρ , any influential round t is $B_{\mathcal{A}}^*$ -free from some agent w.r.t. ρ .*

Proof. Fix realization ρ . Since round t is influential, for some bid profile b and agent i it is (b, ρ) -influential with a strongly influenced agent i . By definition of $b_i^*(\rho, t)$, without loss of generality each bid in b (other than i ’s bid) is at most $b_i^*(\rho, t) \leq B_{\mathcal{A}}^*$. Then $\mathcal{A}(b, \rho, t) \neq i$, and round t is (b, ρ) -secured from agent i .

Suppose round t is not $B_{\mathcal{A}}^*$ -free from agent i w.r.t ρ . Then there exists a bid profile b' in which each bid (other than i ’s bid) is at least $B_{\mathcal{A}}^*$ such that $\mathcal{A}(b', \rho, t) = i$. To derive a contradiction, let us transform b to b' by adjusting first the bid of agent i and then bids of agents $j \neq i$ one agent at a time. Initially agent i is not chosen in round t , and after the last step of this transformation agent i is chosen. Thus it is chosen at some step, say when we adjust the bid of agent i or some agent $j \neq i$. This *transfer of impression* to agent i cannot happen when bid of agent i is adjusted from b_i to b'_i (since round t is (b, ρ) -secured from i), and it cannot happen when bid of player $j \neq i$ is adjusted from b_j to $b'_j \geq b_j$ (this is because, the transfer to i cannot happen from j because of pointwise-monotonicity and the transfer to i cannot happen from $l \neq j$ because of IIA). This is a contradiction. \square

Let T be the time horizon. Assume $v_{\max} \geq 2B_{\mathcal{A}}^*$. Let $N(\rho)$ be the number of influential rounds w.r.t realization ρ . Let $N_i(\rho)$ be the number of influential rounds w.r.t. realization ρ that are $B_{\mathcal{A}}^*$ -free from agent i w.r.t. ρ . Then N and the N_i ’s are random variables in the probability space induced by the clicks. By Lemma C.2 we have that $\sum_i N_i(\rho)$ is at least the number of *influential rounds*. As in Section 4, let $\vec{\mu}_0$ be the vector of CTRs in which all CTRs are $\frac{1}{2}$, and let $\mathbb{E}_0[\cdot]$ denote expectation w.r.t. $\vec{\mu}_0$.

Fix a constant $\beta > 0$ to be specified later. If $\mathbb{E}_0[N] \geq \beta k T^{2/3}$ then $\mathbb{E}_0[N_i] \geq \beta T^{2/3}$ for some agent i , so the allocation incurs expected regret $R(T; v_{\max}) \geq \Omega(v_{\max} T^{2/3})$ on any problem instance \mathcal{J}_j , $j \neq i$. (In this problem instance, CTRs given by $\vec{\mu}_0$, the bid of agent j is v_{\max} , and all other bids are $v_{\max}/2$.) Now suppose $\mathbb{E}_0[N] \leq \beta k T^{2/3}$. Then the desired regret bound follows by an argument very similar to the one in the last paragraph of the proof of Theorem 4.1.

Appendix D: Universally truthful randomized mechanisms

Consider randomized mechanisms that are *universally truthful*, i.e. truthful for each realization of the internal random seed. For mechanisms that randomize over exploration-separated deterministic mechanisms, we obtain the same lower bounds as in Theorems 4.1 and Theorem 4.3.

Theorem D.1. *Consider the MAB mechanism design problem. Let \mathcal{D} distribution over exploration-separated deterministic allocation rules. Then*

$$\mathbb{E}_{\mathcal{A} \in \mathcal{D}} [R_{\mathcal{A}}(T; v_{\max})] = \Omega(v_{\max} k^{1/3} T^{2/3}).$$

Proof Sketch. Recall that in the proof of Theorem 4.1 we define a family \mathcal{F} of $2k$ problem instances, and show that if \mathcal{A} is an exploration-separated deterministic allocation rule, then on one of these instances its regret is “high”. In fact, we can extend this analysis to show that the regret is “high”, that is at least $R^* = \Omega(v_{\max} k^{1/3} T^{2/3})$, on an instance $\mathcal{I} \in \mathcal{F}$ chosen uniformly at random from \mathcal{F} ; here regret is in expectation over the choice of \mathcal{I} .¹⁵ Once this is proved, it follows that regret is $R^*/2$ for any *distribution* over such \mathcal{A} , in expectation over both the choice of \mathcal{A} and the choice of \mathcal{I} . Thus there exists a single (deterministic) instance \mathcal{I} such that $\mathbb{E}_{\mathcal{A} \in \mathcal{D}} [R_{\mathcal{A}, \mathcal{I}}(T)] \geq R^*/2$. \square

Theorem 4.3 extends similarly.

Appendix E: Randomized allocations and adversarial clicks

In this section we discuss randomized allocations and the version of the MAB mechanism design problem when clicks are generated adversarially, termed the *adversarial MAB problem*. In this version, the objective is to optimize the worst-case regret over all values $v = (v_1, \dots, v_k)$ such that $v_i \in [0, v_{\max}]$ for each i , and all realizations ρ :

$$R(T; v; \rho) = \left[\max_i v_i \sum_{t=1}^T \rho_i(t) \right] - \sum_{t=1}^T \sum_{i=1}^k v_i \rho_i(t) \mathbb{E} [\mathcal{A}_i(v; \rho; t)] \quad (\text{E.1})$$

$$R(T; v_{\max}) = \max \{ R(T; v; \rho) : \text{all realizations } \rho, \text{ all } v \text{ such that } v_i \in [0, v_{\max}] \text{ for each } i \}.$$

The first term in (E.1) is the social welfare from the best time-invariant allocation, the second term is the social welfare generated by \mathcal{A} .

Let us make a few definitions related to truthfulness. Recall that a mechanism is called *weakly truthful* if for each realization, it is truthful in expectation over its random seed. A randomized allocation is *pointwise monotone* if for each realization and each bid profile, increasing the bid of any one agent does not decrease the probability of this agent being allocated in any given round. For a set S of rounds and a function $\sigma : S \rightarrow \{\text{agents}\}$, an allocation is (S, σ) -*separated* if (i) it coincides with σ on S , (ii) the clicks from the rounds not in S are discarded (not reported to the algorithm). An allocation is *strongly separated* if before round 1, without looking at the bids, it randomly chooses a set S of rounds and a function $\sigma : S \rightarrow \{\text{agents}\}$, and then runs a pointwise monotone (S, σ) -separated allocation. Note that the choice of S and σ is independent of the clicks, by definition.

We show that for any (randomized) strongly separated allocation rule \mathcal{A} there exists a payment rule which results in a mechanism that is weakly truthful and normalized. Then we consider PSIM [8, 25], a randomized MAB algorithm from the literature, and show that it is pointwise monotone and strongly separated. When interpreted as an allocation rule, there algorithm has strong regret guarantees for the *adversarial MAB mechanism design problem*, where the clicks are chosen by an *oblivious adversary*. Specifically, PSIM obtains regret $R(T, v_{\max}) = O(v_{\max} k^{1/3} (\log k)^{1/3} T^{2/3})$.

We start with the structural result.

¹⁵This extension requires but minor modifications to the proof of Theorem 4.1. For instance, for the case $k \geq 3$ we argue that first, if $\mathbb{E}_0[N] > R$ then $\mathbb{E}_0[N_i] \leq \frac{2}{k} E_0[N]$ for at least $\frac{k}{2}$ agents i (and so on), and if $\mathbb{E}_0[N] \leq R$ then (omitting some details) there are $\Omega(k)$ good agents i such that $\mathbb{E}_0[N_i] \leq 2R/k$ (and so on).

Lemma E.1. *Consider the MAB mechanism design problem. Let \mathcal{A} be a (randomized) strongly separated allocation rule. Then there exists a payment rule \mathcal{P} such that the resulting mechanism $(\mathcal{A}, \mathcal{P})$ is normalized and weakly truthful.*

Proof. Throughout the proof, let us fix a realization ρ , time horizon T , bid vector b , and agent i . We will consider the payment of agent i . We will vary the bid of agent i on the interval $[0, b_i]$; the bids b_{-i} of all other agents always stay the same.

Let $c_i(x)$ be the number of clicks received by agent i given that her bid is x . Then by (the appropriate version of) Theorem 3.1 the payment of agent i must be $\mathcal{P}_i(b)$ such that

$$\mathbb{E}_{\mathcal{A}}[\mathcal{P}_i(b)] = \mathbb{E}_{\mathcal{A}} \left[b_i c_i(b_i) - \int_{x=0}^{b_i} c_i(x) dx \right], \quad (\text{E.2})$$

where the expectation is taken over the internal randomness in the algorithm.

Recall that initially \mathcal{A} randomly selects, without looking at the bids, a set S of rounds and a function $\sigma : S \rightarrow \{\text{agents}\}$, and then runs some pointwise monotone (S, σ) -separated allocation $\mathcal{A}^{(S, \sigma)}$. In what follows, let us fix S and σ , and denote $\mathcal{A}^* = \mathcal{A}^{(S, \sigma)}$. We will refer to the rounds in S as *exploration rounds*, and to the rounds not in S as *exploitation rounds*. Let $\gamma_i^*(x, t)$ be the probability that algorithm \mathcal{A}^* allocates agent i in round t given that agent i bids x . Note that for fixed value of internal random seed of \mathcal{A}^* this probability can only depend on the clicks observed in exploration rounds, which are known to the mechanism. Therefore, abstracting away the computational issues, we can assume that it is known to the mechanism. Define the payment rule as follows: in each exploitation round t in which agent i is chosen and clicked, charge

$$\mathcal{P}_i^*(b, t) = b_i - \frac{1}{\gamma_i^*(b_i, t)} \int_0^{b_i} \gamma_i^*(x, t) dx. \quad (\text{E.3})$$

Then the total payment assigned to agent i is

$$\mathcal{P}_i^*(b) = \sum_{t \notin S} \rho_i(t) \mathcal{A}_i^*(b; \rho; t) \mathcal{P}_i^*(b, t). \quad (\text{E.4})$$

Since allocation \mathcal{A}^* is pointwise monotone, the probability $\gamma_i^*(x, t)$ is non-decreasing in x . Therefore $\mathcal{P}_i^*(b, t) \in [0, b_i]$ for each round t . It follows that the mechanism is normalized (for any realization of the random seed of allocation \mathcal{A}).

It remains to check that the payment rule (E.3) results in (E.2). Let $c_i^*(x)$ be the number of clicks allocated to agent i by allocation \mathcal{A}^* given that her bid is x . Let $c_i^{\text{expl}}(x)$ be the corresponding number of clicks in exploitation rounds only. Since \mathcal{A}^* is (S, σ) -separated, we have

$$\mathbb{E}[c_i^*(x) - c_i^{\text{expl}}(x)] = \sum_{t \in S} \rho_{\sigma(t)}(t) = \text{const}(x). \quad (\text{E.5})$$

Taking expectations in (E.4) over the random seed of \mathcal{A}_S and using (E.5), we obtain

$$\begin{aligned} \mathbb{E}[\mathcal{P}_i^*(b)] &= \sum_{t \notin S} \rho_i(t) \gamma_i^*(b_i, t) \mathcal{P}_i^*(b, t) \\ &= \sum_{t \notin S} \rho_i(t) \left[b_i \gamma_i^*(b_i, t) - \int_0^{b_i} \gamma_i^*(x, t) dx \right] \\ &= b_i \left[\sum_{t \notin S} \rho_i(t) \gamma_i^*(b_i, t) \right] - \int_0^{b_i} \left[\sum_{t \notin S} \rho_i(t) \gamma_i^*(x, t) \right] dx \\ &= b_i \mathbb{E}[c_i^{\text{expl}}(b_i)] - \int_0^{b_i} \mathbb{E}[c_i^{\text{expl}}(x)] dx \\ &= \mathbb{E} \left[b_i c_i^*(b_i) - \int_0^{b_i} c_i^*(x) dx \right]. \end{aligned}$$

Finally, taking expectations over the choice of S and σ , we obtain (E.2). \square

E.1 Algorithm PSIM is strongly separated

In this subsection, we consider PSIM [8, 25], an algorithm for the adversarial MAB problem. We interpret this algorithm as an allocation rule, and observe that it is strongly separated.

As usual, k denotes the number of agents; let $[k]$ denote the set of agents.

Input: Time horizon T , bid vector b . Let $v_{\max} = \max_i b_i$.

Output: For each round $t \leq T$, a distribution on $[k]$.

1. Divide the time horizon into P phases of T/P consecutive rounds each.
2. From rounds of each phase p , pick without replacement k rounds at random (called the *exploration rounds*) and assign them randomly to k arms. Let S denote the set of all exploration rounds (of all phases). Let $f : S \rightarrow [k]$ be the function which tells which arm is assigned to an exploration round in S . The rounds in $[T] \setminus S$ are called the exploitation rounds.
3. Let $w_i(0) = 1$ for all $i \in [k]$.
4. For each phase $p = 1, 2, \dots, P$
 - (a) For each round t in phase p
 - i. If $t \in S$ and $f(t) = i$, then define the distribution $\gamma(b; t; S, f)$ such that $\gamma_i(b; t; S, f) = 1$. Pick an agent according to this distribution (equivalently, pick agent i), observe the click $\rho_i(t)$, and update $w_i(p)$ multiplicatively,

$$w_i(p) = w_i(p-1) \cdot (1 + \epsilon)^{\rho_i(t)b_i/v_{\max}}.$$

- ii. If $t \notin S$, then define the distribution $\gamma(b; t; S, f)$ such that $\gamma_i(b; t; S, f) = \frac{w_i(p-1)}{\sum_j w_j(p-1)}$. Pick an agent according to $\gamma(b; t; S, f)$, observe the feedback, and discard the feedback.

If we pick the values $\epsilon = (k \log k / T)^{1/3}$ and $P = (\log k)^{1/3} (T/k)^{2/3}$, then the regret of PSIM is bounded by $\mathcal{O}((k \log k)^{1/3} T^{2/3} v_{\max})$ against any oblivious adversary (see [8, 25]).

We next prove that PSIM is strongly-separated.

It is clear from the structure of PSIM above that it chooses a set S of exploration rounds and a function $f : S \rightarrow [k]$ in the beginning without looking at the bids and then runs an (S, f) -separated allocation. We need to prove that the (S, f) -separated allocation is pointwise monotone. For this we need prove that the probability $\gamma_i(b; t; S, f)$ is monotone in the bid of agent i , where $\gamma_i(b; t; S, f)$ denotes the probability of picking agent i in round t when bids are b given the choice of S and f . If $t \in S$, the $\gamma_i(b; t; S, f)$ is independent of bids, and hence is monotone in b_i . Let $t \notin S$ and t is a round in phase p . Let us denote by $f^{-1}(i, p)$ the (unique) exploration round in phase p assigned to agent i . We then have

$$\gamma_i(b; t; S, f) = (1 + \epsilon)^{\frac{b_i}{v_{\max}} \sum_{q=1}^{p-1} \rho_i(f^{-1}(i, q))} \Bigg/ \sum_j (1 + \epsilon)^{\frac{b_j}{v_{\max}} \sum_{q=1}^{p-1} \rho_j(f^{-1}(j, q))}.$$

We split the denominator into the term for agent i and all other terms. It is then not hard to see that this is a non-decreasing function of b_i .

We state the above results in the form of the following corollary.

Corollary E.2. *There exists a weakly truthful normalized mechanism for the adversarial MAB problem (against oblivious adversary) whose regret grows as $\mathcal{O}((k \log k)^{1/3} \cdot T^{2/3} \cdot v_{\max})$.*

Appendix F: Truthfulness in expectation over CTRs

We consider the stochastic MAB mechanism design problem under a more relaxed notion of truthfulness: truthfulness *in expectation*, where for each vector of CTRs the expectation is taken over clicks (and the internal randomness in the mechanism, if the latter is not deterministic). We show that any allocation \mathcal{A}^* that is monotone in expectation,¹⁶ can be converted to a mechanism that is truthful in expectation and monotone in expectation, with minor changes and a very minor increase in regret. Furthermore, we show that there exist MAB allocations that are monotone in expectation whose regret matches the optimal upper bounds for MAB *algorithms*. The conclusion is that in order to obtain any non-trivial lower bounds on regret and (essentially) any non-trivial structural results, one needs to assume that a mechanism is ex-post normalized, at least in some approximate sense.

The main result of this section is that for any allocation \mathcal{A}^* that is monotone in expectation, any time horizon T , and any parameter $\gamma \in (0, 1)$ there exists a mechanism $(\mathcal{A}, \mathcal{P})$ such that the mechanism is truthful in expectation and normalized in expectation, and allocation \mathcal{A} initially makes a random choice between \mathcal{A}^* and some other allocation, choosing \mathcal{A}^* with probability at least γ . We call such allocation \mathcal{A} a γ -approximation of \mathcal{A}^* . Clearly, on any problem instance we have $R_{\mathcal{A}}(T) \leq \gamma R_{\mathcal{A}^*}(T) + (1 - \gamma)T$. The extra additive factor of $(1 - \gamma)T$ is not significant if e.g. $\gamma = 1 - \frac{1}{T}$. The problem with this mechanism is that it is not ex-post normalized; moreover, in some realizations payments may be very large in absolute value.

Theorem F.1. *Consider the stochastic MAB mechanism design problem with k agents and a fixed time horizon T . For each $\gamma \in (0, 1)$ and each allocation rule \mathcal{A}^* that is monotone in expectation, there exists a mechanism $(\mathcal{A}, \mathcal{P})$ such that \mathcal{A} is a γ -approximation of \mathcal{A}^* , and the mechanism is truthful in expectation and normalized in expectation.*

Remark. Payment rule \mathcal{P} is well-defined as a mapping from histories to numbers. We do not make any claims on the efficient computability thereof.

For the sake of completeness, we provide a concrete algorithm which one could plug into Theorem F.1 and obtain improved (and in fact, best possible) regret guarantees.

Proposition F.2. *Consider the stochastic MAB mechanism design problem with k agents and a fixed time horizon T . There exists an allocation rule \mathcal{A} that is monotone in expectation, whose regret is $R(T; v_{\max}) = O(v_{\max} \sqrt{kT \log T})$ in the worst case, and $R_{\delta}(T; v_{\max}) = O(v_{\max} \frac{k}{\delta} \log T)$ on the δ -gap instances.*

Proof Sketch. For simplicity, assume $v_{\max} = 1$. Let $r_0 = \sqrt{8 \log(T)/T}$. Consider the following simple allocation. Initially, each agent is *active*. In each phase, play each active agent once, in a round-robin fashion. After the phase, (permanently) de-activate each agent whose *sample product* (sample average times the bid) is more than r_0 below that of some other active agent. This completes the description of the allocation.

This allocation is based on a well-known (perhaps folklore) MAB algorithm. The regret bounds are proved along the lines of those in [6]. The crucial observations are that with a very high probability the optimal agent is never de-activated, and that each sub-optimal agent i is played at most $O(\Delta_i^{-2} \log T)$ times, where Δ_i is the difference between her product (CTR times the bid) and the maximal one.

The allocation is monotone in expectation because increasing the bid of a given agent cannot cause this agent to be de-activated later. \square

¹⁶Monotonicity in expectation is defined in an obvious way: an allocation is *monotone in expectation* if for each agent i and fixed bid profile b_{-i} , the corresponding expected click-allocation is a non-decreasing function of b_i ; here the expectation is taken over the clicks and possibly the allocation's random seed.

F.1 Proof of Theorem F.1

Let $\mathcal{A}_{\text{expl}}$ be the allocation rule where in each round an agent is chosen independently and uniformly at random. Allocation \mathcal{A} is defined as follows: use \mathcal{A}^* with probability γ ; otherwise use $\mathcal{A}_{\text{expl}}$. Fix an instance (b, μ) of the stochastic MAB mechanism design problem, where $b = (b_1, \dots, b_k)$ and $\mu = (\mu_1, \dots, \mu_k)$ are vectors of bids and CTRs, respectively. Let $C_i = C_i(b_i; b_{-i})$ be the expected number of clicks for agent i under the original allocation \mathcal{A}^* . Then by Myerson [33] the expected payment of agent i must be

$$\mathcal{P}_i^M = \gamma \left[b_i C_i(b_i; b_{-i}) - \int_0^{b_i} C_i(x; b_{-i}) dx \right]. \quad (\text{F.1})$$

The key idea is to treat the expected payment as a multivariate polynomial over μ_1, \dots, μ_k . It is essential (given the way we define \mathcal{P}) to show that this polynomial has degree $\leq T$.

Claim F.3. \mathcal{P}_i^M is a polynomial of degree $\leq T$ in variables μ_1, \dots, μ_k .

Proof. Fix the bid profile. Let X_t be allocation of algorithm \mathcal{A}^* . Let $\text{poly}(T)$ be the set of all polynomials over μ_1, \dots, μ_k of degree at most T . Consider a fixed history $h = (x_1, y_1; \dots; x_T, y_T)$, and let h^t be the corresponding history up to (and including) round t . Then

$$\mathbb{P}[h] = \prod_{t=1}^T \Pr[X_t = x_t \mid h^{t-1}] \mu_{x_t}^{y_t} (1 - \mu_{x_t})^{1-y_t} \in \text{poly}(T) \quad (\text{F.2})$$

$$C_i(b_i; b_{-i}) = \sum_{h \in \mathcal{H}} \mathbb{P}[h] \# \text{clicks}_i(h) \in \text{poly}(T). \quad (\text{F.3})$$

Therefore $\mathcal{P}_i^M \in \text{poly}(T)$, since one can take an integral in (F.1) separately over the coefficient of each monomial of $C_i(x; b_{-i})$. \square

Fix time horizon T . For a given run of an allocation rule, the *history* is defined as $h = (x_1, y_1; \dots; x_T, y_T)$, where x_t is the allocation in round t , and $y_t \in \{0, 1\}$ is the corresponding click. Let \mathcal{H} be the set of all possible histories.

Our payment rule \mathcal{P} is a deterministic function of history. For each agent i , we define the payment $\mathcal{P}_i = \mathcal{P}_i(h)$ for each history h such that $E_h[\mathcal{P}_i(h)] = \mathcal{P}_i^M$ for any choice of CTRs, and hence $E_h[\mathcal{P}_i(h)] \equiv \mathcal{P}_i^M$, where \equiv denotes an equality between polynomials over μ_1, \dots, μ_k .

Fix the bid vector and fix agent i . We define the payment \mathcal{P}_i as follows. Charge nothing if allocation \mathcal{A}^* is used. If allocation $\mathcal{A}_{\text{expl}}$ is used, charge *per monomial*. Specifically, let $\text{mono}(T)$ be the set of all monomials over μ_1, \dots, μ_k of degree at most T . For each monomial $Q \in \text{mono}(T)$ we define a subset of *relevant histories* $\mathcal{H}_i(Q) \subset \mathcal{H}$. (We defer the definition till later in the proof.) For a given history $h \in \mathcal{H}$ we charge a (possibly negative) amount

$$\mathcal{P}_i(h) = \frac{1}{1-\gamma} \sum_{Q \in \text{mono}(T): h \in \mathcal{H}_i(Q)} k^{\deg(Q)} \mathcal{P}_i^M(Q), \quad (\text{F.4})$$

where $\deg(Q)$ is the degree of Q , and $\mathcal{P}_i^M(Q)$ is the coefficient of Q in \mathcal{P}_i^M . Let \mathbb{P}_{expl} be the distribution on histories induced by $\mathcal{A}_{\text{expl}}$. Then the expected payment is

$$E_h[\mathcal{P}_i(h)] = \sum_{Q \in \text{mono}(T)} k^{\deg(Q)} \mathbb{P}_{\text{expl}}[\mathcal{H}_i(Q)] \mathcal{P}_i^M(Q).$$

Therefore in order to guarantee that $E_h[\mathcal{P}_i(h)] \equiv \mathcal{P}_i^M$ it suffices to choose $\mathcal{H}_i(Q)$ for each Q so that

$$k^{\deg(Q)} \mathbb{P}_{\text{expl}}[\mathcal{H}_i(Q)] \equiv Q. \quad (\text{F.5})$$

Consider a monomial $Q = \mu_1^{\alpha_1} \dots \mu_k^{\alpha_k}$. Let $\mathcal{H}_i(Q)$ consist of all histories such that first agent 1 is played α_1 times in a row, and clicked every time, then agent 2 is played α_2 times in a row, and clicked every time, and so on till agent k . In the remaining $T - \deg(Q)$ rounds, any agent can be chosen, and any outcome (click or no click) can be received. It is clear that (F.5) holds.